

# Distributed adaptive steplength stochastic approximation schemes for Cartesian stochastic variational inequality problems

Farzad Yousefian, Angelia Nedić, and Uday V. Shanbhag\*

January 10, 2013

## Abstract

Motivated by problems arising in decentralized control problems and non-cooperative Nash games, we consider a class of strongly monotone Cartesian variational inequality (VI) problems, where the mappings either contain expectations or their evaluations are corrupted by error. Such complications are captured under the umbrella of Cartesian stochastic variational inequality problems and we consider solving such problems via stochastic approximation (SA) schemes. Specifically, we propose a scheme wherein the steplength sequence is derived by a rule that depends on problem parameters such as monotonicity and Lipschitz constants. The proposed scheme is seen to produce sequences that are guaranteed to converge almost surely to the unique solution of the problem. To cope with networked multi-agent generalizations, we provide requirements under which independently chosen steplength rules still possess desirable almost-sure convergence properties. In the second part of this paper, we consider a regime where Lipschitz constants on the map are either unavailable or difficult to derive. Here, we present a local randomization technique that allows for deriving an approximation of the original mapping, which is then shown to be Lipschitz continuous with a prescribed constant. Using this technique, we introduce a locally randomized SA algorithm and provide almost sure convergence theory for the resulting sequence of iterates to an approximate solution of the original variational inequality problem. Finally, the paper concludes with some preliminary numerical results on a stochastic rate allocation problem and a stochastic Nash-Cournot game.

## 1 Introduction

Multi-agent system-theoretic problems can collectively capture a range of problems arising from decentralized control problems and noncooperative games. In static regimes, where agent problems are convex and agent feasibility sets are uncoupled, the associated solutions of such problems are given by the solution of a suitably defined Cartesian variational inequality problem. Our interest lies in settings where the mapping arising in such problems is strongly monotone and one of the following hold: (i) Either the mapping contains expectations whose analytical form is unavailable; or (ii) The evaluation of such a mapping is corrupted by error. In either case, the appropriate problem of interest is given by a stochastic variational inequality problem  $\text{VI}(X, F)$  that requires determining an  $x^* \in X$  such that

$$(x - x^*)^T F(x^*) \geq 0 \quad \text{for all } x \in X, \quad (1)$$

where

$$F(x) \triangleq \begin{pmatrix} \mathbb{E}[\Phi_1(x, \xi)] \\ \vdots \\ \mathbb{E}[\Phi_N(x, \xi)] \end{pmatrix}, \quad (2)$$

---

\*The first two authors are with the Department of Industrial and Enterprise Systems Engineering, University of Illinois, Urbana, IL 61801, USA, while the last author is with the Department of Industrial and Manufacturing Engineering, Pennsylvania State University, University Park, PA 16802, USA. They are contactable at {yousefi1, angelia}@illinois.edu and udaybag@psu.edu. Nedić and Shanbhag gratefully acknowledge the support of the NSF through the award NSF CMMI 0948905 ARRA. Additionally, Nedić has been funded by NSF award CMMI-0742538 and Shanbhag has been supported by NSF award CMMI-1246887.

$\Phi_i : \mathcal{D}_i \times \mathbb{R}^d \rightarrow \mathbb{R}^{n_i}$ ,  $\mathcal{D}_i \subseteq \mathbb{R}^{n_i}$ ,  $X$  is a closed and convex set,  $\mathcal{D}_i$  is an open set in  $\mathbb{R}^{n_i}$  and  $\sum_{i=1}^N n_i = n$ . Furthermore,  $\xi : \Omega \rightarrow \mathbb{R}^d$  is a random variable, where  $\Omega$  denotes the associated sample space and  $\mathbb{E}[\cdot]$  denotes the expectation with respect to  $\xi$ .

Variational inequality problems assume relevance in capturing the solution sets of convex optimization and equilibrium problems [11]. Their Cartesian specializations arise from specifying the set  $X$  as a Cartesian product, i.e.,  $X \triangleq \prod_{i=1}^N X_i$ . Such problems arise in the modeling of multi-agent decision-making problems such as rate allocation problems in communication networks [18, 31, 35], noncooperative Nash games in communication networks [1, 2, 39], competitive interactions in cognitive radio networks [20, 29, 30, 38], and strategic behavior in power markets [16, 17, 32]. Our interest lies in regimes complicated by uncertainty, which could arise as a result of agents facing expectation-based objectives that do not have tractable analytical forms. Naturally, the Cartesian stochastic variational inequality problem framework represents an expansive model for capturing a range of such problems.

Two broad avenues exist for solving such a class of problems. Of these, the first approach, referred to as the sample-average approximation (SAA) method. In adopting this approach, one uses a set of  $M$  samples  $\{\xi_1, \dots, \xi_M\}$  and considers the sample-average problem where an expected mapping  $\mathbb{E}[\Phi(x, \xi)]$  is replaced by the sample-average  $\sum_{j=1}^M \Phi(x, \xi^j)/M$ . The resulting problem is deterministic and its solution provides an estimator for the solution of the true problem. The asymptotic behavior of these estimators has been studied extensively in the context of stochastic optimization and variational problems [23, 33]. The other approach, referred to as stochastic approximation, also has a long tradition. First proposed by Robbins and Monro [28] for root-finding problems and by Ermoliev for stochastic programs [8–10], significant effort has been applied towards theoretical and algorithmic examination of such schemes (cf. [4, 21, 34]). Yet, there has been markedly little on the application of such techniques to solution of stochastic variational inequalities, exceptions being [14, 19]. Standard stochastic approximation schemes provide little guidance regarding the choice of a steplength sequence, denoted by  $\{\gamma_k\}$ , apart from requiring that the sequence satisfies

$$\sum_{k=0}^{\infty} \gamma_k = \infty \quad \text{and} \quad \sum_{k=0}^{\infty} \gamma_k^2 < \infty.$$

The behavior of stochastic approximation schemes is closely tied to the choice of steplength sequences. Generally, there have been two avenues traversed in choosing steplengths: (i) *Deterministic steplength sequences*: Spall [34, Ch. 4, pg. 113] considered diverse choices of the form  $\gamma_k = \frac{\beta}{(k+1+a)^\alpha}$ , where  $\beta > 0$ ,  $0 < \alpha \leq 1$ , and  $a \geq 0$  is a stability constant. In related work in the context of approximate dynamic programming, Powell [27] examined several deterministic update rules. However, much of these results are not provided with convergence theory. (ii) *Stochastic steplength sequences*: An alternative to a deterministic rule is a stochastic scheme that updates steplengths based on observed data. Of note is recent work by George et al. [12] where an adaptive stepsize rule is proposed that minimizes the mean squared error. In a similar vein, Cicek et al. [7] develop an adaptive Kiefer-Wolfowitz SA algorithm and derive general upper bounds on its mean-squared error.

Before proceeding, we note the relationship of the present work to three specific references. In [19], Cartesian stochastic variational inequality problems with Lipschitzian mappings were considered with a focus towards integrating Tikhonov and prox-based regularization techniques with standard stochastic gradient methods. However, the steplength sequences were “non-adaptive” since the choices did not adapt to problem parameters. Two problem-specific adaptive rules were developed in our earlier work on stochastic convex programming. Additionally, local smoothing techniques were examined for addressing the lack of smoothness. Of these, the first, referred to as the *recursive steplength* SA scheme, forms the inspiration for a generalization pursued in the current work. Finally, in [40], we extended this recursive rule to accommodate stochastic variational inequality problems. Note that the qualifier “adaptive” implies that the steplength rule adapts to problem parameters such as Lipschitz constant, monotonicity constant and the diameter of the set. In this paper, our goal lies in developing a distributed adaptive stochastic approximation scheme (DASA) that can accommodate networked multi-agent implementations and cope with non-Lipschitzian mappings. More specifically, the main contributions of this paper are as follows:

(i) **DASA schemes for Lipschitzian CSVIs**: We begin with a simple extension of the adaptive steplength rule presented in [41] to the variational regime under a Lipschitzian requirement on the map. Yet, implementing this rule in a centralized regime is challenging and this motivates the need for distributed coun-

terparts that can be employed on Cartesian problems. Such a distributed rule is developed and produces sequences of iterates that are guaranteed to converge to the solution in almost-sure sense.

**(ii) DASA schemes for non-Lipschitzian CSVIs:** Our second goal lies in addressing the absence or unavailability of a Lipschitz constant by leveraging locally randomized smoothing techniques, again inspired by our efforts to solve nonsmooth stochastic optimization problems [41]. In this part of the paper, we generalize this natively centralized scheme for optimization problems to a distributed version that can cope with Cartesian stochastic variational inequality problems.

The remainder of this paper is organized as follows. In Section 2, we provide a canonical formulation for the problem of interest and motivate this formulation through two sets of examples. An adaptive steplength SA scheme for stochastic variational inequality problems with Lipschitzian mappings and its distributed generalization are provided in Section 3. By leveraging a locally randomized smoothing technique, in Section 4, we extend these schemes to a regime where Lipschitzian assumptions do not hold. Finally, the paper concludes with some preliminary numerics in Section 5.

**Notation:** Throughout this paper, a vector  $x$  is assumed to be a column vector. We write  $x^T$  to denote the transpose of a vector  $x$ ,  $\|x\|$  to denote the Euclidean vector norm, i.e.,  $\|x\| = \sqrt{x^T x}$ ,  $\|x\|_1$  to denote the 1-norm, i.e.,  $\|x\|_1 = \sum_{i=1}^n |x_i|$  for  $x \in \mathbb{R}^n$ , and  $\|x\|_\infty$  to denote the infinity vector norm, i.e.,  $\|x\|_\infty = \max_{i=1, \dots, n} |x_i|$  for  $x \in \mathbb{R}^n$ . We use  $\Pi_X(x)$  to denote the Euclidean projection of a vector  $x$  on a set  $X$ , i.e.,  $\|x - \Pi_X(x)\| = \min_{y \in X} \|x - y\|$ . For a convex function  $f$  with domain  $\text{dom} f$ , a vector  $g$  is a *subgradient* of  $\bar{x} \in \text{dom} f$  if  $f(\bar{x}) + g^T(x - \bar{x}) \leq f(x)$  holds for all  $x \in \text{dom} f$ . The set of all subgradients of  $f$  at  $\bar{x}$  is denoted by  $\partial f(\bar{x})$ . We write *a.s.* as the abbreviation for “almost surely”. We use  $\text{Prob}(A)$  to denote the probability of an event  $A$  and  $\mathbb{E}[z]$  to denote the expectation of a random variable  $z$ . The **Matlab** notation  $(u_1; u_2; u_3)$  refers to a column vector with components  $u_1$ ,  $u_2$  and  $u_3$ , respectively.

## 2 Formulation and source problems

In Section 2.1, we formulate the Cartesian stochastic variational inequality (CSVI) problem and outline the stochastic approximation algorithmic framework. A motivation for studying CSVIs is provided through two examples in Section 2.2, while a review of the main assumptions is given in Section 2.3.

### 2.1 Problem formulation and algorithm outline

Given a set  $X \subseteq \mathbb{R}^n$  and a mapping  $F : X \rightarrow \mathbb{R}^n$ , the variational inequality problem, denoted by  $\text{VI}(X, F)$ , requires determining a vector  $x^* \in X$  such that  $(x - x^*)^T F(x^*) \geq 0$  holds for all  $x \in X$ . When the underlying set  $X$  is given by a Cartesian product, as articulated by the definition  $X \triangleq \prod_{i=1}^N X_i$ , where  $X_i \subseteq \mathbb{R}^{n_i}$ , then the associated variational inequality is qualified as a *Cartesian* variational inequality problem. Now suppose that  $x^* = (x_1^*; x_2^*; \dots; x_N^*) \in X$  satisfies the following system of inequalities:

$$(x_i - x_i^*)^T \mathbb{E}[\Phi_i(x^*, \xi_i)] \geq 0 \quad \text{for all } x_i \in X_i \text{ and all } i = 1, \dots, N, \quad (3)$$

where  $\xi_i : \Omega_i \rightarrow R^{d_i}$  is a random vector with some probability distribution for  $i = 1, \dots, N$ . Naturally, problem (3) may be reduced to  $\text{VI}(X, F)$  by noting that  $F$  may be defined as in (2), where  $n = \sum_{i=1}^N n_i$  and  $F : X \rightarrow \mathbb{R}^n$ . Then,  $\text{VI}(X, F)$  is a stochastic variational inequality problem on the Cartesian product of the sets  $X_i$  with a solution  $x^* = (x_1^*; x_2^*; \dots; x_N^*)$ .

Much of the interest in this paper pertains to the development of stochastic approximation schemes for  $\text{VI}(X, F)$  when the components the map  $F$  is defined by (2). For such a problem, we consider the following distributed stochastic approximation scheme:

$$\begin{aligned} x_{k+1,i} &= \Pi_{X_i}(x_{k,i} - \gamma_{k,i}(F_i(x_k) + w_{k,i})), \\ w_{k,i} &\triangleq \Phi_i(x_k, \xi_{k,i}) - F_i(x_k), \end{aligned} \quad (4)$$

for all  $k \geq 0$  and  $i = 1, \dots, N$ , where  $F_i(x) \triangleq \mathbb{E}[\Phi_i(x, \xi_i)]$  for  $i = 1, \dots, N$ ,  $\gamma_{k,i} > 0$  is the stepsize for the  $i$ th index at iteration  $k$ ,  $x_{k,i}$  denotes the solution for the  $i$ -th index at iteration  $k$ , and  $x_k = (x_{k,1}; x_{k,2}; \dots; x_{k,N})$ . Moreover,  $x_0 \in X$  is a random initial vector independent of any other random variables in the scheme and such that  $\mathbb{E}[\|x_0\|^2] < \infty$ .

## 2.2 Motivating examples

We consider two problems that can be addressed by Cartesian stochastic variational inequality framework.

**Example 1** (Networked stochastic Nash-Cournot game). A classical example of a Nash game is a networked Nash-Cournot game [15, 24]. Suppose a collection of  $N$  firms compete over a network of  $M$  nodes wherein the production and sales for firm  $i$  at node  $j$  are denoted by  $g_{ij}$  and  $s_{ij}$ , respectively. Suppose firm  $i$ 's cost of production at node  $j$  is denoted by the uncertain cost function  $c_{ij}(g_{ij}, \xi)$ . Furthermore, goods sold by firm  $i$  at node  $j$  fetch a random revenue defined by  $p_j(\bar{s}_j, \xi)s_{ij}$  where  $p_j(\bar{s}_j, \xi)$  denotes the uncertain sales price at node  $j$  and  $\bar{s}_j = \sum_{i=1}^N s_{ij}$  denotes the aggregate sales at node  $j$ . Finally, firm  $i$ 's production at node  $j$  is capacitated by  $\text{cap}_{ij}$  and its optimization problem is given by the following<sup>1</sup>:

$$\begin{aligned} & \text{minimize} && \mathbb{E}[f_i(x, \xi)] \\ & \text{subject to} && x_i \in X_i, \end{aligned}$$

where  $x = (x_1; \dots; x_N)$  with  $x_i = (g_i; s_i)$ ,  $g_i = (g_{i1}; \dots; g_{iM})$ ,  $s_i = (s_{i1}; \dots; s_{iM})$ , and

$$\begin{aligned} f_i(x, \xi) &\triangleq \sum_{j=1}^M (c_{ij}(g_{ij}, \xi) - p_j(\bar{s}_j, \xi)s_{ij}), \\ X_i &\triangleq \left\{ (g_i, s_i) \mid \sum_{j=1}^M g_{ij} = \sum_{j=1}^M s_{ij}, \quad g_{ij}, s_{ij} \geq 0, \quad g_{ij} \leq \text{cap}_{ij}, \quad j = 1, \dots, M \right\}. \quad \square \end{aligned}$$

Under the validity of the interchange between the expectation and the derivative operator, the resulting equilibrium conditions of this stochastic Nash-Cournot game are compactly captured by the variational inequality  $\text{VI}(X, F)$  where  $X \triangleq \prod_{i=1}^N X_i$  and  $F(x) = (F_1(x); \dots; F_N(x))$  with  $F_i(x) = \mathbb{E}[\nabla_{x_i} f_i(x, \xi)]$ .

**Example 2** (Stochastic composite minimization problem). Consider a generalized min-max optimization problem given by

$$\text{minimize} \quad \Psi(\psi_1(x), \dots, \psi_m(x)) \tag{5}$$

$$\text{subject to} \quad x \in X \triangleq \prod_{i=1}^N X_i, \tag{6}$$

where  $\Psi(u_1, \dots, u_m)$  is defined as

$$\Psi(u_1, \dots, u_m) \triangleq \max_{y \in \mathcal{Y}} \left\{ \sum_{i=1}^m u_i^T (A_i y + b_i) - \beta(y) \right\},$$

while  $\psi_i(x) \triangleq \mathbb{E}[\phi_i(x, \xi)]$ ,  $\nabla_{x_j} \psi_i(x) = \mathbb{E}[H_{ji}(x, \xi)]$  for  $i = 1, \dots, m$ , and  $\beta(y)$  is a Lipschitz continuous convex function of  $y$ .  $\square$

Under the assumption that the derivative and the expectation operator can be interchanged, it can be seen that the solution to this optimization problem can be obtained by solving a Cartesian stochastic variational inequality problem  $\text{VI}(X \times \mathcal{Y}, F)$  where

$$F(x, y) \triangleq \begin{pmatrix} \sum_{i=1}^m \nabla_{x_1} \psi_i(x) (A_i y + b_i) \\ \vdots \\ \sum_{i=1}^m \nabla_{x_N} \psi_i(x) (A_i y + b_i) \\ - \sum_{i=1}^m A_i^T \psi_i(x) + \nabla_y \beta(y) \end{pmatrix} = \begin{pmatrix} \mathbb{E}[\sum_{i=1}^m H_{1i}(x, \xi) (A_i y + b_i)] \\ \vdots \\ \mathbb{E}[\sum_{i=1}^m H_{Ni}(x, \xi) (A_i y + b_i)] \\ \mathbb{E}[-\sum_{i=1}^m A_i^T \phi_i(x, \xi) + \nabla_y \beta(y)] \end{pmatrix}.$$

Note that the specification that  $\psi_i(x)$  and its Jacobian are expectation-valued may be a consequence of not having access to noise-free evaluations of either object. In particular, one only has access to evaluations  $\phi_i(x, \xi)$  and Jacobian evaluations given by  $H_{ji}(x, \xi) = \nabla_{x_j} \phi_i(x, \xi)$ .

<sup>1</sup>Note that the transportation costs are assumed to be zero.

## 2.3 Assumptions

Our interest lies in the development of distributed stochastic approximation schemes for Cartesian stochastic variational inequality problems as espoused by (4) and the associated global convergence theory in regimes where the mappings are single-valued mappings that are not necessarily Lipschitz continuous. We let

$$X = \prod_{i=1}^N X_i,$$

and make the following assumptions on the set  $X$  and the mapping  $F$ .

**Assumption 1.** *Assume the following:*

- (a) *The set  $X_i \subseteq \mathbb{R}^{n_i}$  is closed and convex for  $i = 1, \dots, N$ .*
- (b) *The mapping  $F(x)$  is a single-valued Lipschitz continuous over the set  $X$  with a constant  $L$ .*
- (c) *The mapping  $F(x)$  is strongly monotone with a constant  $\eta > 0$ :*

$$(F(x) - F(y))^T(x - y) \geq \eta \|x - y\|^2 \quad \text{for all } x, y \in X.$$

Since  $F$  is strongly monotone, the existence and uniqueness of the solution to  $\text{VI}(X, F)$  is guaranteed by Theorem 2.3.3 of [11]. We let  $x^*$  denote the solution of  $\text{VI}(X, F)$ .

Regarding the method in (4), we let  $\mathcal{F}_k$  denote the history of the method up to time  $k$ , i.e.,  $\mathcal{F}_k = \{x_0, \xi_0, \xi_1, \dots, \xi_{k-1}\}$  for  $k \geq 1$  and  $\mathcal{F}_0 = \{x_0\}$ , where  $\xi_k = (\xi_{k,1}; \xi_{k,2}; \dots; \xi_{k,N})$ . In terms of this definition, we note that

$$\mathbb{E}[w_{k,i} \mid \mathcal{F}_k] = \mathbb{E}[\Phi_i(x_k, \xi_{k,i}) \mid \mathcal{F}_k] - F_i(x_k) = 0 \quad \text{for all } k \geq 0 \text{ and all } i.$$

We impose some further conditions on the stochastic errors  $w_{k,i}$  of the algorithm, as follows.

**Assumption 2.** *The errors  $w_k = (w_{k,1}; w_{k,2}; \dots; w_{k,N})$  are such that for some (deterministic)  $\nu > 0$ ,*

$$\mathbb{E}[\|w_k\|^2 \mid \mathcal{F}_k] \leq \nu^2 \quad \text{a.s. for all } k \geq 0.$$

We use the following Lemma in establishing the convergence of method (4) and its extensions. This result may be found in [26] (cf. Lemma 10, page 49).

**Lemma 1.** *Let  $\{v_k\}$  be a sequence of nonnegative random variables, where  $\mathbb{E}[v_0] < \infty$ , and let  $\{\alpha_k\}$  and  $\{\mu_k\}$  be deterministic scalar sequences such that:*

$$\mathbb{E}[v_{k+1} \mid v_0, \dots, v_k] \leq (1 - \alpha_k)v_k + \mu_k \quad \text{a.s. for all } k \geq 0,$$

$$0 \leq \alpha_k \leq 1, \quad \mu_k \geq 0, \quad \text{for all } k \geq 0,$$

$$\sum_{k=0}^{\infty} \alpha_k = \infty, \quad \sum_{k=0}^{\infty} \mu_k < \infty, \quad \lim_{k \rightarrow \infty} \frac{\mu_k}{\alpha_k} = 0.$$

*Then,  $v_k \rightarrow 0$  almost surely,  $\lim_{k \rightarrow \infty} \mathbb{E}[v_k] = 0$ , and for any  $\epsilon > 0$  and for all  $k > 0$ ,*

$$\text{Prob}(\{v_j \leq \epsilon \text{ for all } j \geq k\}) \geq 1 - \frac{1}{\epsilon} \left( \mathbb{E}[v_k] + \sum_{i=k}^{\infty} \beta_i \right).$$

## 3 Distributed adaptive SA schemes for Lipschitzian mappings

In this section, we restrict our attention to settings where the mapping  $F(x)$  is a single-valued Lipschitzian map. In Section 3.1, we begin by developing an adaptive steplength rule for deriving steplength sequences from problem parameters such as monotonicity constant, Lipschitz constant etc., where the qualifier adaptive implies that the steplength choices “adapt” or are “self-tuned” to problem parameters. Unfortunately, in distributed regimes, such a rule requires prescription by a central coordinator, a relatively challenging task in multi-agent regimes. This motivates the development of a distributed counterpart of the aforementioned adaptive rule and provide convergence theory for such a generalization in Section 3.2.

### 3.1 An adaptive steplength SA (ASA) scheme

Stochastic approximation algorithms require stepsize sequences to be square summable but not summable. These algorithms provide little advice regarding the choice of such sequences. One of the most common choices has been the harmonic steplength rule which takes the form of  $\gamma_k = \frac{\theta}{k}$  where  $\theta > 0$  is a constant. Although, this choice guarantees almost-sure convergence, it does not leverage problem parameters. Numerically, it has been observed that such choices can perform quite poorly in practice. Motivated by this shortcoming, we present a steplength scheme for a centralized variant of algorithm (4):

$$\begin{aligned} x_{k+1} &= \Pi_X(x_k - \gamma_k(F(x_k) + w_k)), \\ w_k &\triangleq \Phi(x_k, \xi_k) - F(x_k), \end{aligned} \quad (7)$$

for  $k \geq 0$ . The proposed scheme derives a rule for updating steplength sequences that adapts to problem parameters while guaranteeing almost-sure convergence of  $x_k$  to the unique solution of  $\text{VI}(X, F)$ .

A key challenge in practical implementations of stochastic approximation lies in choosing an appropriate diminishing steplength sequence  $\{\gamma_k\}$ . In [41], we developed a rule for selecting such a sequence in a convex stochastic optimization regime by leveraging three parameters: (i) Lipschitz constant of the gradients; (ii) strong convexity constant; and (ii) diameter of the set  $X$ . Along similar directions, such a rule is constructed for strongly monotone stochastic variational inequality problems and the results in this subsection bear significant similarity to those presented in [41] with some key distinctions. First, these results are presented for strongly monotone stochastic variational inequality problems and second, co-coercivity of the mappings is not assumed, leading to a tighter requirement on the choice of steplengths.

**Lemma 2.** *Consider algorithm (7), and let Assumptions 1 and 2 hold. Then, the following relation holds almost surely for all  $k \geq 0$ :*

$$\mathbb{E}[\|x_{k+1} - x^*\|^2 \mid \mathcal{F}_k] \leq (1 - 2\eta\gamma_k + L^2\gamma_k^2)\|x_k - x^*\|^2 + \gamma_k^2\nu^2. \quad (8)$$

*Proof.* By the definition of algorithm (7) and the non-expansiveness property of the projection operator, we have for all  $k \geq 0$ ,

$$\begin{aligned} \|x_{k+1} - x^*\|^2 &= \|\Pi_X(x_k - \gamma_k(F(x_k) + w_k)) - \Pi_X(x^* - \gamma_k F(x^*))\|^2 \\ &\leq \|x_k - x^* - \gamma_k(F(x_k) + w_k - F(x^*))\|^2. \end{aligned}$$

Taking expectations conditioned on the past, and by employing  $\mathbb{E}[w_k \mid \mathcal{F}_k] = 0$ , we have

$$\begin{aligned} \mathbb{E}[\|x_{k+1} - x^*\|^2 \mid \mathcal{F}_k] &\leq \|x_k - x^*\|^2 + \gamma_k^2\|F(x_k) - F(x^*)\|^2 + \gamma_k^2\mathbb{E}[\|w_k\|^2 \mid \mathcal{F}_k] \\ &\quad - 2\gamma_k(x_k - x^*)^T(F(x_k) - F(x^*)) \\ &\leq (1 - 2\eta\gamma_k + \gamma_k^2L^2)\|x_k - x^*\|^2 + \gamma_k^2\nu^2, \end{aligned}$$

where the second inequality is a consequence of the strong monotonicity and Lipschitz continuity of  $F(x)$  over  $X$  as well as the boundedness of  $\mathbb{E}[\|w_k\|^2 \mid \mathcal{F}_k]$ .  $\square$

The upper bound (8) can be used to construct an adaptive stepsize rule. Note that inequality (8) holds for any  $\gamma_k > 0$ . When the stepsize is further restricted so that  $0 < \gamma_k \leq \frac{\eta}{L^2}$ , we have

$$1 - \gamma_k(2\eta - \gamma_kL^2) \leq 1 - \eta\gamma_k.$$

Thus, for  $0 < \gamma_k \leq \frac{\eta}{L^2}$  and by taking expectations, inequality (8) reduces to

$$\mathbb{E}[\|x_{k+1} - x^*\|^2] \leq (1 - \eta\gamma_k)\mathbb{E}[\|x_k - x^*\|^2] + \gamma_k^2\nu^2 \quad \text{for all } k \geq 0. \quad (9)$$

We begin by viewing  $\mathbb{E}[\|x_{k+1} - x^*\|^2]$  as an error  $e_{k+1}$  arising from employing the stepsize sequence  $\gamma_0, \gamma_1, \dots, \gamma_k$ . Furthermore, the worst case error arises when (9) holds as an equality and satisfies the following recursive relation:

$$e_{k+1}(\gamma_0, \dots, \gamma_k) = (1 - \eta\gamma_k)e_k(\gamma_0, \dots, \gamma_{k-1}) + \gamma_k^2\nu^2.$$

Motivated by this relationship, our interest lies in examining whether the stepsizes  $\gamma_0, \gamma_1, \dots, \gamma_k$  can be chosen so as to minimize the error  $e_k$ . Our goal lies in constructing a stepsize scheme that allows for claiming the almost sure convergence of the sequence  $\{x_k\}$  produced by algorithm (7) to the unique solution  $x^*$  of  $\text{VI}(X, F)$ . We formalize this approach by defining real-valued error functions  $e_k(\gamma_0, \dots, \gamma_{k-1})$  as follows:

$$e_k(\gamma_0, \dots, \gamma_{k-1}) \triangleq (1 - \eta\gamma_{k-1})e_{k-1}(\gamma_0, \dots, \gamma_{k-2}) + \gamma_{k-1}^2\nu^2 \quad \text{for } k \geq 1, \quad (10)$$

where  $e_0$  is a positive scalar,  $\eta$  is the strong monotonicity constant and  $\nu^2$  is an upper bound for the second moments of the error norms  $\|w_k\|$ . We consider a choice of  $\{\gamma_0, \gamma_1, \dots, \gamma_{k-1}\}$  based on minimizing an upper bound on the mean-squared error, namely  $e(\gamma_0, \gamma_1, \dots, \gamma_{k-1})$ , as captured by the following optimization problem:

$$\begin{aligned} & \text{minimize} && e_k(\gamma_0, \dots, \gamma_{k-1}) \\ & \text{subject to} && 0 < \gamma_j \leq \frac{\eta}{L^2} \text{ for all } j = 0, \dots, k-1. \end{aligned}$$

To ensure convergence in an almost-sure sense, the sequence  $\{\gamma_k\}$  needs to satisfy  $\sum_{j=0}^{\infty} \gamma_j = \infty$  and  $\sum_{j=0}^{\infty} \gamma_j^2 < \infty$ . As the next two propositions show, these can indeed be achieved. In fact, the error  $e_{k+1}$  at the next iteration can also be minimized by selecting  $\gamma_k$  as a function of only the most recent stepsize  $\gamma_{k-1}$ . In what follows, we consider the sequence  $\{\gamma_k^*\}$  given by

$$\gamma_0^* = \frac{\eta}{2\nu^2} e_0 \quad (11)$$

$$\gamma_k^* = \gamma_{k-1}^* \left(1 - \frac{\eta}{2}\gamma_{k-1}^*\right) \quad \text{for all } k \geq 1. \quad (12)$$

We provide a result showing that the stepsizes  $\gamma_i$ ,  $i = 0, \dots, k-1$ , minimize  $e_k$  over  $(0, \frac{\eta}{L^2}]^k$ , where  $L$  is the Lipschitz constant associated with  $F(x)$  over  $X$ .

**Proposition 1** (An adaptive steplength SA (ASA) scheme). *Let the error function  $e_k(\gamma_0, \dots, \gamma_{k-1})$  be defined as in (10), where  $e_0 \geq 0$  is such that  $\nu \geq L\sqrt{\frac{e_0}{2}}$ , where  $L$  is the Lipschitz constant of  $F$ . Let the sequence  $\{\gamma_k^*\}$  be given by (11)–(12). Then, the following hold:*

- (a) *For all  $k \geq 0$ , the error  $e_k$  satisfies  $e_k(\gamma_0^*, \dots, \gamma_{k-1}^*) = \frac{2\nu^2}{\eta} \gamma_k^*$ .*
- (b) *For each  $k \geq 1$ , the vector  $(\gamma_0^*, \gamma_1^*, \dots, \gamma_{k-1}^*)$  is the minimizer of the function  $e_k(\gamma_0, \dots, \gamma_{k-1})$  over the set*

$$\mathbb{G}_k \triangleq \left\{ \alpha \in \mathbb{R}^k : 0 < \alpha_j \leq \frac{\eta}{L^2} \text{ for } j = 1, \dots, k \right\}.$$

*More precisely, for any  $k \geq 1$  and any  $(\gamma_0, \dots, \gamma_{k-1}) \in \mathbb{G}_k$ , we have*

$$e_k(\gamma_0, \dots, \gamma_{k-1}) - e_k(\gamma_0^*, \dots, \gamma_{k-1}^*) \geq \nu^2(\gamma_{k-1} - \gamma_{k-1}^*)^2.$$

The almost-sure convergence of the produced sequence holds for a family of steplength rules, as captured by the following result.

**Proposition 2** (Almost-sure convergence of ASA scheme). *Let Assumptions 1 and 2 hold. Assume that the stepsize sequence  $\{\gamma_k\}$  is generated by the following adaptive scheme:*

$$\gamma_k = \gamma_{k-1}(1 - c\gamma_{k-1}) \quad \text{for all } k \geq 1,$$

*where  $c > 0$  is a scalar and the initial stepsize is such that  $0 < \gamma_0 < \frac{1}{c}$ . Then, the sequence  $\{x_k\}$  generated by algorithm (7) converges almost surely to a random point that belongs to the optimal set.*

The proofs of Propositions 1 and 2 are omitted, as they follow from a more general results for a distributed SA method, as discussed in the next subsection.



### 3.2 A distributed adaptive steplength SA (DASA) scheme

Unfortunately, in multi-agent regimes, the implementation of the stepsize rule (11)-(12) requires a central coordinator who can prescribe and enforce such rules. In this section, we extend the centralized rule to accommodate a multi-agent setting wherein each agent chooses its own update rule, given the global knowledge of problem parameters. In such a regime, given that the set  $X$  is the Cartesian product of closed and convex sets  $X_1, \dots, X_N$ , our interest lies in developing steplength update rules in the context of method (4) where the  $i$ -th agent chooses its steplength, denoted by  $\gamma_{k,i}$ , as per

$$\gamma_{k,i} = \gamma_{k-1,i}(1 - c_i \gamma_{k-1,i}),$$

with  $c_i > 0$  being a constant associated with agent  $i$  mapping  $F_i(x)$ , while the initial stepsize  $\gamma_{0,i}$  is suitably selected. The following assumption imposes requirements on the stepsizes  $\gamma_{k,i}$  in (4).

**Assumption 3.** Assume that the following hold:

(3a) The stepsize sequences  $\{\gamma_{k,i}\}$ ,  $i = 1, \dots, N$ , are such that  $\gamma_{k,i} > 0$  for all  $k$  and  $i$ , with  $\sum_{k=0}^{\infty} \gamma_{k,i} = \infty$  and  $\sum_{k=0}^{\infty} \gamma_{k,i}^2 < \infty$  for all  $i$ .

(3b) If  $\{\delta_k\}$  and  $\{\Gamma_k\}$  are positive sequences such that  $\delta_k \leq \min_{1 \leq i \leq N} \gamma_{k,i}$  and  $\Gamma_k \geq \max_{1 \leq i \leq N} \gamma_{k,i}$  for all  $k \geq 0$ , then

$$\frac{\Gamma_k - \delta_k}{\delta_k} \leq \beta \quad \text{for all } k \geq 0,$$

where  $\beta$  is a scalar satisfying  $0 \leq \beta < \frac{\eta}{L}$ .

**Remark:** Assumption (3a) is a standard requirement on steplength sequences while Assumption (3b) provides an additional condition on the discrepancy between the stepsize values  $\gamma_{k,i}$  at each iteration  $k$ . This condition is satisfied, for instance, when  $\gamma_{k,1} = \dots = \gamma_{k,N}$ , in which case  $\beta = 0$ .

When deriving an adaptive rule, we use Lemma 1 and a distributed generalization of Lemma 2, which is given below.

**Lemma 3.** Consider algorithm (4). Let Assumptions 1 and 2 hold.

(a) The following relation holds almost surely for all  $k \geq 0$ :

$$\mathbb{E}[\|x_{k+1} - x^*\|^2 \mid \mathcal{F}_k] \leq (1 - 2(\eta + L)\delta_k + 2L\Gamma_k + L^2\Gamma_k^2)\|x_k - x^*\|^2 + \Gamma_k^2\nu^2,$$

where  $\{\delta_k\}$  and  $\{\Gamma_k\}$  are positive sequences such that  $\delta_k \leq \min_{1 \leq i \leq N} \gamma_{k,i}$  and  $\Gamma_k \geq \max_{1 \leq i \leq N} \gamma_{k,i}$  for all  $k$ .

(b) If Assumption (3b) holds, then the following relation is valid for all  $k \geq 0$ :

$$\mathbb{E}[\|x_{k+1} - x^*\|^2] \leq (1 - 2(\eta - \beta L)\delta_k + (1 + \beta)^2 L^2 \delta_k^2) \mathbb{E}[\|x_k - x^*\|^2] + (1 + \beta)^2 \delta_k^2 \nu^2.$$

*Proof.* (a) From the properties of the projection operator, we know that a vector  $x^*$  solves  $\text{VI}(X, F)$  problem if and only if  $x^*$  satisfies  $x^* = \Pi_X(x^* - \gamma F(x^*))$  for any  $\gamma > 0$ . By the definition of algorithm (4) and the non-expansiveness property of the projection operator, we have for all  $k \geq 0$  and all  $i$ ,

$$\begin{aligned} \|x_{k+1,i} - x_i^*\|^2 &= \|\Pi_{X_i}(x_{k,i} - \gamma_{k,i}(F_i(x_k) + w_{k,i})) - \Pi_{X_i}(x_i^* - \gamma_{k,i}F_i(x^*))\|^2 \\ &\leq \|x_{k,i} - x_i^* - \gamma_{k,i}(F_i(x_k) + w_{k,i} - F_i(x^*))\|^2. \end{aligned}$$

Taking the expectation conditioned on the past, and using  $\mathbb{E}[w_{k,i} \mid \mathcal{F}_k] = 0$ , we have

$$\begin{aligned} \mathbb{E}[\|x_{k+1,i} - x_i^*\|^2 \mid \mathcal{F}_k] &\leq \|x_{k,i} - x_i^*\|^2 + \gamma_{k,i}^2 \|F_i(x_k) - F_i(x^*)\|^2 + \gamma_{k,i}^2 \mathbb{E}[\|w_{k,i}\|^2 \mid \mathcal{F}_k] \\ &\quad - 2\gamma_{k,i}(x_{k,i} - x_i^*)^T (F_i(x_k) - F_i(x^*)). \end{aligned}$$



Now, by summing the preceding relations over  $i$ , we have

$$\begin{aligned} \mathbb{E}[\|x_{k+1} - x^*\|^2 \mid \mathcal{F}_k] &\leq \|x_k - x^*\|^2 + \sum_{i=1}^N \gamma_{k,i}^2 \|F_i(x_k) - F_i(x^*)\|^2 + \sum_{i=1}^N \gamma_{k,i}^2 \mathbb{E}[\|w_{k,i}\|^2 \mid \mathcal{F}_k] \\ &\quad - 2 \sum_{i=1}^N \gamma_{k,i} (x_{k,i} - x_i^*)^T (F_i(x_k) - F_i(x^*)). \end{aligned}$$

Using  $\gamma_{k,i} \leq \Gamma_k$  and Assumption 2, we can see that  $\sum_{i=1}^N \gamma_{k,i}^2 \mathbb{E}[\|w_{k,i}\|^2 \mid \mathcal{F}_k] \leq \Gamma_k^2 \nu^2$  almost surely for all  $k \geq 0$ . Thus, from the preceding relation, we have

$$\begin{aligned} \mathbb{E}[\|x_{k+1} - x^*\|^2 \mid \mathcal{F}_k] &\leq \|x_k - x^*\|^2 + \underbrace{\sum_{i=1}^N \gamma_{k,i}^2 \|F_i(x_k) - F_i(x^*)\|^2}_{\text{Term 1}} + \Gamma_k^2 \nu^2 \\ &\quad - \underbrace{2 \sum_{i=1}^N \gamma_{k,i} (x_{k,i} - x_i^*)^T (F_i(x_k) - F_i(x^*))}_{\text{Term 2}}. \end{aligned} \quad (13)$$

Next, we estimate Term 1 and Term 2 in (13). By using the definition of  $\Gamma_k$  and by leveraging the Lipschitzian property of mapping  $F$ , we obtain

$$\text{Term 1} \leq \Gamma_k^2 \|F(x_k) - F(x^*)\|^2 \leq \Gamma_k^2 L^2 \|x_k - x^*\|^2. \quad (14)$$

By adding and subtracting  $-2 \sum_{i=1}^N \delta_k (x_{k,i} - x_i^*)^T (F_i(x_k) - F_i(x^*))$  from Term 2, and using  $\sum_{i=1}^N (x_{k,i} - x_i^*)^T (F_i(x_k) - F_i(x^*)) = (x_k - x^*)^T (F(x_k) - F(x^*))$ , we further obtain

$$\text{Term 2} \leq -2\delta_k (x_k - x^*)^T (F(x_k) - F(x^*)) - 2 \sum_{i=1}^N (\gamma_{k,i} - \delta_k) (x_{k,i} - x_i^*)^T (F_i(x_k) - F_i(x^*)).$$

By Cauchy-Schwartz inequality, the preceding relation yields

$$\begin{aligned} \text{Term 2} &\leq -2\delta_k (x_k - x^*)^T (F(x_k) - F(x^*)) + 2(\gamma_{k,i} - \delta_k) \sum_{i=1}^N \|x_{k,i} - x_i^*\| \|F_i(x_k) - F_i(x^*)\| \\ &\leq -2\delta_k (x_k - x^*)^T (F(x_k) - F(x^*)) + 2(\Gamma_k - \delta_k) \|x_k - x^*\| \|F(x_k) - F(x^*)\|, \end{aligned}$$

where in the last relation, we use the definition of  $\Gamma_k$  and Hölder's inequality. Invoking strong monotonicity of the mapping  $F$  for bounding the first term and by utilizing the Lipschitzian property of the second term of the preceding relation, we have

$$\text{Term 2} \leq -2\eta\delta_k \|x_k - x^*\|^2 + 2(\Gamma_k - \delta_k)L \|x_k - x^*\|^2. \quad (15)$$

The desired inequality is obtained by combining relations (13), (14), and (15).

(b) Assumption 3b implies that  $\Gamma_k \leq (1 + \beta)\delta_k$ . Combining this observation with the result of part (a), we obtain almost surely for all  $k \geq 0$ ,

$$\mathbb{E}[\|x_{k+1} - x^*\|^2 \mid \mathcal{F}_k] \leq (1 - 2(\eta - \beta L)\delta_k + (1 + \beta)^2 L^2 \delta_k^2) \|x_k - x^*\|^2 + (1 + \beta)^2 \delta_k^2 \nu^2.$$

Taking expectations in the preceding inequality, we obtain the desired relation.  $\square$

The following proposition proves the almost-sure convergence of the distributed SA scheme when the steplength sequences satisfy the bounds prescribed by Assumption 3b.

**Proposition 3** (Almost-sure convergence of distributed SA scheme). *Let Assumptions 1, 2, and 3 hold. Then, the sequence  $\{x_k\}$  generated by algorithm (4) converges almost surely to the unique solution of  $VI(X, F)$ .*

*Proof.* Consider the relation of Lemma 3(a). For this relation, we show that the conditions of Lemma 1 are satisfied, which will allow us to claim the almost-sure convergence of  $x_k$  to  $x^*$ . Let us define  $v_k \triangleq \|x_k - x^*\|^2$ , and

$$\alpha_k \triangleq 2(\eta - \beta L)\delta_k - L^2\delta_k^2(1 + \beta)^2, \quad \mu_k \triangleq (1 + \beta)^2\delta_k^2\nu^2. \quad (16)$$

Next, we show that  $0 \leq \alpha_k \leq 1$  for  $k$  sufficiently large. Since  $\gamma_{k,i}$  tends to zero for all  $i = 1, \dots, N$ , we may conclude that  $\delta_k$  goes to zero as  $k$  grows. In turn, as  $\delta_k$  goes to zero, for  $k$  large enough, say  $k \geq k_1$ , we have

$$1 - \frac{(1 + \beta)^2 L^2 \delta_k}{2(\eta - \beta L)} > 0.$$

By Assumption 3b we have  $\beta < \frac{\eta}{L}$ , which implies  $\eta - \beta L > 0$ . Thus, we have  $\alpha_k \geq 0$  for  $k \geq k_1$ . Also, for  $k$  large enough, say  $k \geq k_2$ , we have  $\alpha_k \leq 1$  since  $\delta_k \rightarrow 0$ . Therefore, when  $k \geq \max\{k_1, k_2\}$  we have  $0 \leq \alpha_k \leq 1$ . Obviously,  $v_k, \mu_k \geq 0$ .

From Assumption 3b we have  $\delta_k \leq \gamma_k \leq (1 + \beta)\delta_k$  for all  $k$ . Using these relations and the conditions on  $\gamma_{k,i}$  given in Assumption 3a, we can show that  $\sum_{k=0}^{\infty} \delta_k = \infty$  and  $\sum_{k=0}^{\infty} \delta_k^2 < \infty$ . Furthermore, from the preceding properties of the sequence  $\{\delta_k\}$ , and the definitions of  $\alpha_k$  and  $\mu_k$  in (16), we can see that  $\sum_{k=0}^{\infty} \alpha_k = \infty$  and  $\sum_{k=0}^{\infty} \mu_k < \infty$ . Finally, by the definitions of  $\alpha_k$  and  $\mu_k$  we have

$$\lim_{k \rightarrow \infty} \frac{\mu_k}{\alpha_k} = \lim_{k \rightarrow \infty} \left( \frac{(1 + \beta)^2 \delta_k \nu^2}{2(\eta - \beta L) \left(1 - \frac{(1 + \beta)^2 L^2 \delta_k}{2(\eta - \beta L)}\right)} \right) = \frac{(1 + \beta)^2 (\lim_{k \rightarrow \infty} \delta_k) \nu^2}{2(\eta - \beta L) \left(1 - \frac{(1 + \beta)^2 L^2 (\lim_{k \rightarrow \infty} \delta_k)}{2(\eta - \beta L)}\right)},$$

implying that  $\lim_{k \rightarrow \infty} \frac{\mu_k}{\alpha_k} = 0$  since  $\delta_k \rightarrow 0$ . Hence, all conditions of Lemma 1 are satisfied and we may conclude that  $\|x_k - x^*\|^2 \rightarrow 0$  almost surely.  $\square$

Proposition 3 states that under specified assumptions on the set  $X$  and mapping  $F$ , the stochastic errors  $w_k$ , and the stepsizes  $\gamma_{k,i}$ , the distributed SA scheme is guaranteed to converge to the unique solution of  $\text{VI}(X, F)$  almost surely. Our goal in the remainder of this section lies in providing a stepsize rule that aims to minimize a suitably defined error function of the algorithm, while satisfying Assumption 3. To begin our analysis, we consider the result of Lemma 3b for all  $k \geq 0$ :

$$\mathbb{E}[\|x_{k+1} - x^*\|^2] \leq (1 - 2(\eta - \beta L)\delta_k + (1 + \beta)^2 L^2 \delta_k^2) \mathbb{E}[\|x_k - x^*\|^2] + (1 + \beta)^2 \delta_k^2 \nu^2, \quad (17)$$

where  $\delta_k \leq \min_{1 \leq i \leq N} \gamma_{k,i}$ . When the stepsizes  $\gamma_{k,i}$  are further restricted so that  $0 < \delta_k \leq \frac{\eta - \beta L}{(1 + \beta)^2 L^2}$ , we have

$$1 - 2(\eta - \beta L)\delta_k + (1 + \beta)^2 L^2 \delta_k^2 \leq 1 - (\eta - \beta L)\delta_k.$$

Thus, for  $0 < \delta_k \leq \frac{\eta - \beta L}{(1 + \beta)^2 L^2}$ , from inequality (17) we obtain

$$\mathbb{E}[\|x_{k+1} - x^*\|^2] \leq (1 - (\eta - \beta L)\delta_k) \mathbb{E}[\|x_k - x^*\|^2] + (1 + \beta)^2 \delta_k^2 \nu^2 \quad \text{for all } k \geq 0. \quad (18)$$

Similar to the discussion in Section 3.1 in the context of the ASA scheme, let us view the quantity  $\mathbb{E}[\|x_{k+1} - x^*\|^2]$  as an error  $e_{k+1}$  of the method arising from the use of the lower bounds  $\delta_0, \delta_1, \dots, \delta_k$  for the stepsize values  $\gamma_{0,i}, \gamma_{1,i}, \dots, \gamma_{k,i}$ ,  $i = 1, \dots, N$ . Relation (18) gives us an error estimate for algorithm (4) in terms of the lower bounds  $\delta_0, \delta_1, \dots, \delta_k$ . We use this estimate to develop an adaptive stepsize procedure. Consider the case when (18) holds with equality, which is the worst case error. In this case, the error satisfies the following recursive relation:

$$e_{k+1} = (1 - (\eta - \beta L)\delta_k) e_k + (1 + \beta)^2 \nu^2 \delta_k^2 \quad \text{for all } k \geq 0.$$

Let us assume that we want to run the algorithm (4) for a fixed number of iterations, say  $K$ . The preceding relation shows that  $e_K$  depends on the lower bound values up to the  $K$ th iteration. This motivates us to view the lower bounds  $\delta_0, \delta_1, \dots, \delta_{K-1}$  as decision variables that can be used to minimize the corresponding upper bound on the mean-squared error of the algorithm up to iteration  $K$ . Thus, the variables are  $\delta_0, \delta_1, \dots, \delta_{K-1}$  and the objective function is the error function  $e_K(\delta_0, \delta_1, \dots, \delta_{K-1})$ . We proceed to derive a rule for generating lower bounds  $\delta_0, \delta_1, \dots, \delta_K$  by minimizing the error  $e_{K+1}$ . Importantly, it turns out that  $\delta_K$  is a

function of only the most recent bound  $\delta_{K-1}$ . We define the real-valued error function  $e_k(\delta_0, \delta_1, \dots, \delta_{k-1})$  by considering an equality in (18):

$$e_{k+1}(\delta_0, \dots, \delta_k) \triangleq (1 - (\eta - \beta L)\delta_k)e_k(\delta_0, \dots, \delta_{k-1}) + (1 + \beta^2)\nu^2\delta_k^2 \quad \text{for all } k \geq 0, \quad (19)$$

where  $e_0$  is a positive scalar,  $\{\delta_k\}$  is a sequence of positive scalars such that  $0 < \delta_k \leq \frac{\eta - \beta L}{(1 + \beta)^2 L^2}$ ,  $L$  is the Lipschitz constant of the mapping  $F$ ,  $\eta$  is the strong monotonicity parameter of  $F$ , and  $\nu^2$  is the upper bound for the second moment of the error norms  $\|w_k\|$  (cf. Assumption 2).

Now let us consider the stepsize sequence  $\{\delta_k^*\}$  given by

$$\delta_0^* = \frac{\eta - \beta L}{2(1 + \beta)^2 \nu^2} e_0 \quad (20)$$

$$\delta_k^* = \delta_{k-1}^* \left( 1 - \left( \frac{\eta - \beta L}{2} \right) \delta_{k-1}^* \right) \quad \text{for all } k \geq 1, \quad (21)$$

where  $e_0$  is the same initial error as for the errors  $e_k$  in (19). In what follows, we often abbreviate  $e_k(\delta_0, \dots, \delta_{k-1})$  by  $e_k$  whenever this is unambiguous. The next proposition shows that the lower bound sequence  $\{\delta_k^*\}$  for  $\gamma_{k,i}$  given by (20)–(21) minimizes the errors  $e_k$  over  $[0, \frac{\eta - \beta L}{(1 + \beta)^2 L^2}]^k$ .

**Proposition 4** (An adaptive lower bound steplength SA scheme). *Let  $e_k(\delta_0, \dots, \delta_{k-1})$  be defined as in (19), where  $e_0$  is a given positive scalar,  $\nu$  is an upper bound defined in Assumption 2,  $\eta$  and  $L$  are the strong monotonicity and Lipschitz constants of the mapping  $F$  respectively and  $\nu$  is chosen such that  $\nu \geq L\sqrt{\frac{e_0}{2}}$ . Let  $\beta$  be a scalar such that  $0 \leq \beta < \frac{\eta}{L}$ , and let the sequence  $\{\delta_k^*\}$  be given by (20)–(21). Then, the following hold:*

- (a) For all  $k \geq 0$ , the error  $e_k$  satisfies  $e_k(\delta_0^*, \dots, \delta_k^*) = \frac{2(1 + \beta)^2 \nu^2}{\eta - \beta L} \delta_k^*$ .
- (b) For any  $k \geq 1$ , the vector  $(\delta_0^*, \delta_1^*, \dots, \delta_{k-1}^*)$  is the minimizer of the function  $e_k(\delta_0, \dots, \delta_{k-1})$  over the set

$$\mathbb{G}_k \triangleq \left\{ \alpha \in \mathbb{R}^k : 0 < \alpha_j \leq \frac{\eta - \beta L}{(1 + \beta)^2 L^2}, j = 1, \dots, k \right\}.$$

More precisely, for any  $k \geq 1$  and any  $(\delta_0, \dots, \delta_{k-1}) \in \mathbb{G}_k$ , we have

$$e_k(\delta_0, \dots, \delta_{k-1}) - e_k(\delta_0^*, \dots, \delta_{k-1}^*) \geq (1 + \beta)^2 \nu^2 (\delta_{k-1} - \delta_{k-1}^*)^2.$$

*Proof.* (a) To show the result, we use induction on  $k$ . Trivially, it holds for  $k = 0$  from (20). Now, suppose that we have  $e_k(\delta_0^*, \dots, \delta_{k-1}^*) = \frac{2(1 + \beta)^2 \nu^2}{\eta - \beta L} \delta_k^*$  for some  $k$ , and consider the case for  $k + 1$ . From the definition of the error  $e_k$  in (19), we have

$$\begin{aligned} e_{k+1}(\delta_0^*, \dots, \delta_k^*) &= (1 - (\eta - \beta L)\delta_k^*)e_k(\delta_0^*, \dots, \delta_{k-1}^*) + (1 + \beta)^2 \nu^2 (\delta_k^*)^2 \\ &= (1 - (\eta - \beta L)\delta_k^*) \frac{2(1 + \beta)^2 \nu^2}{\eta - \beta L} \delta_k^* + (1 + \beta)^2 \nu^2 (\delta_k^*)^2, \end{aligned}$$

where the second equality follows by the inductive hypothesis. Thus,

$$e_{k+1}(\delta_0^*, \dots, \delta_k^*) = \frac{2(1 + \beta)^2 \nu^2}{\eta - \beta L} \delta_k^* \left( 1 - \frac{\eta - \beta L}{2} \delta_k^* \right) = \frac{2(1 + \beta)^2 \nu^2}{\eta - \beta L} \delta_{k+1}^*,$$

where the last equality follows by the definition of  $\delta_{k+1}^*$  in (21). Hence, the result holds for all  $k \geq 0$ .

(b) First we need to show that  $(\delta_0^*, \dots, \delta_{k-1}^*) \in \mathbb{G}_k$ . By our assumption on  $e_0$ , we have  $0 < e_0 \leq \frac{2\nu^2}{L^2}$ , which by the definition of  $\delta_0^*$  in (20) implies that  $0 < \delta_0^* \leq \frac{\eta - \beta L}{(1 + \beta)^2 L^2}$ , i.e.,  $\delta_0^* \in \mathbb{G}_1$ . Using the induction on  $k$ , from relations (20)–(21), it can be shown that  $0 < \delta_k^* < \delta_{k-1}^*$  for all  $k \geq 1$ . Thus,  $(\delta_0^*, \dots, \delta_{k-1}^*) \in \mathbb{G}_k$  for all  $k \geq 1$ . Using the induction on  $k$  again, we now show that the vector  $(\delta_0^*, \delta_1^*, \dots, \delta_{k-1}^*)$  minimizes the

error  $e_k(\delta_0, \dots, \delta_{k-1})$  for all  $k \geq 1$ . From the definition of the error  $e_1$  and the relation  $e_1(\delta_0^*) = \frac{2(1+\beta)^2\nu^2}{\eta-\beta L} \delta_1^*$  shown in part (a), we have

$$e_1(\delta_0) - e_1(\delta_0^*) = (1 - (\eta - \beta L)\delta_0)e_0 + (1 + \beta)^2\nu^2\delta_0^2 - \frac{2(1 + \beta)^2\nu^2}{\eta - \beta L} \delta_1^*.$$

Using  $\delta_1^* = \delta_0^* \left(1 - \frac{\eta - \beta L}{2} \delta_0^*\right)$  (cf. (21)), we obtain

$$e_1(\delta_0) - e_1(\delta_0^*) = (1 - (\eta - \beta L)\delta_0)e_0 + (1 + \beta)^2\nu^2\delta_0^2 - \frac{2(1 + \beta)^2\nu^2}{\eta - \beta L} \delta_0^* + (1 + \beta)^2\nu^2(\delta_0^*)^2.$$

Since  $e_0 = \frac{2(1+\beta)^2\nu^2}{\eta-\beta L} \delta_0^*$  (cf. (20)), it follows that

$$e_1(\delta_0) - e_1(\delta_0^*) = -2(1 + \beta)^2\nu^2\delta_0\delta_0^* + (1 + \beta)^2\nu^2\delta_0^2 + (1 + \beta)^2\nu^2(\delta_0^*)^2 = (1 + \beta)^2\nu^2(\delta_0 - \delta_0^*)^2,$$

showing that the inductive hypothesis holds for  $k = 1$ . Now, suppose that

$$e_k(\delta_0, \dots, \delta_{k-1}) - e_k(\delta_0^*, \dots, \delta_{k-1}^*) \geq (1 + \beta)^2\nu^2(\delta_{k-1} - \delta_{k-1}^*)^2. \quad (22)$$

holds for some  $k$  and for all  $(\delta_0, \dots, \delta_{k-1}) \in \mathbb{G}_k$ . We next show that relation (22) holds for  $k + 1$  and for all  $(\delta_0, \dots, \delta_k) \in \mathbb{G}_{k+1}$ . To simplify the notation, we use  $e_{k+1}^*$  to denote the error  $e_{k+1}$  evaluated at  $(\delta_0^*, \delta_1^*, \dots, \delta_k^*)$ , and  $e_{k+1}$  when evaluating at an arbitrary vector  $(\delta_0, \delta_1, \dots, \delta_k) \in \mathbb{G}_{k+1}$ . Using (19) and part (a), we have

$$e_{k+1} - e_{k+1}^* = (1 - (\eta - \beta L)\delta_k)e_k + (1 + \beta)^2\nu^2\delta_k^2 - \frac{2(1 + \beta)^2\nu^2}{\eta - \beta L} \delta_{k+1}^*.$$

Under the inductive hypothesis, we have  $e_k \geq e_k^*$  (cf. (22)). When  $(\delta_0, \delta_1, \dots, \delta_k) \in \mathbb{G}_k$ , we have  $\delta_k \leq \frac{(\eta - \beta L)}{(1 + \beta)^2 L^2}$ . Next, we show that  $\frac{(\eta - \beta L)}{(1 + \beta)^2 L^2} \leq \frac{1}{\eta - \beta L}$ . By the definition of strong monotonicity and Lipschitzian property, we have  $\eta \leq L$ . Using  $\eta \leq L$  and  $0 \leq \beta \leq \frac{\eta}{L}$  we obtain

$$\begin{aligned} \eta &\leq (1 + \beta)L \Rightarrow \eta - \beta L \leq (1 + \beta)L \\ \Rightarrow (\eta - \beta L)^2 &\leq (1 + \beta)^2 L^2 \Rightarrow \frac{(\eta - \beta L)}{(1 + \beta)^2 L^2} \leq \frac{1}{\eta - \beta L}. \end{aligned}$$

This implies that for  $(\delta_0, \delta_1, \dots, \delta_k) \in \mathbb{G}_k$ , we have  $\delta_k \leq \frac{1}{\eta - \beta L}$  or equivalently  $1 - (\eta - \beta L)\delta_k \geq 0$ . Using this, the relation  $e_k^* = \frac{2(1+\beta)^2\nu^2}{\eta-\beta L} \delta_k^*$  of part (a), and the definition of  $\delta_{k+1}^*$ , we obtain

$$\begin{aligned} e_{k+1} - e_{k+1}^* &\geq (1 - (\eta - \beta L)\delta_k) \frac{2(1 + \beta)^2\nu^2}{\eta - \beta L} \delta_k^* + (1 + \beta)^2\nu^2\delta_k^2 - \frac{2(1 + \beta)^2\nu^2}{\eta - \beta L} \delta_k^* \left(1 - \frac{\eta - \beta L}{2} \delta_k^*\right) \\ &= (1 + \beta)^2\nu^2(\delta_k - \delta_k^*)^2. \end{aligned}$$

Hence, we have  $e_k - e_k^* \geq (1 + \beta)^2\nu^2(\delta_{k-1} - \delta_{k-1}^*)^2$  for all  $k \geq 1$  and all  $(\delta_0, \dots, \delta_{k-1}) \in \mathbb{G}_k$ .  $\square$

**Remark:** From Proposition 4, the minimizer  $(\delta_0^*, \dots, \delta_{k-1}^*)$  of  $e_k$  over  $\mathbb{G}_k$  is unique up to a scaling by a factor  $\rho \in (0, 1)$ . Specifically, the solution  $(\delta_0^*, \dots, \delta_{k-1}^*)$  is obtained for an initial error  $e_0 \geq 0$  satisfying  $\nu \geq L \sqrt{\frac{e_0}{2}}$ , where  $e_0$  can be chosen to be arbitrarily large by scaling  $\nu$  appropriately. Suppose that in the definition of the sequence  $\{\delta_k^*\}$ ,  $\rho e_0$  is employed instead of  $e_0$  for some  $\rho \in (0, 1)$ . Then it can be seen (by following the proof) that, for the resulting sequence, Proposition 4 would still hold.  $\square$

We have just provided an analysis in terms of a lower bound sequence  $\{\delta_k\}$ . We may conduct a similar analysis for an upper bound sequence  $\{\Gamma_k\}$ . In particular, from Lemma 3a we have

$$\mathbb{E}[\|x_{k+1} - x^*\|^2] \leq (1 - 2(\eta + L)\delta_k + 2L\Gamma_k + L^2\Gamma_k^2)\mathbb{E}[\|x_k - x^*\|^2] + \Gamma_k^2\nu^2 \quad \text{for all } k \geq 0.$$

When  $\frac{\Gamma_k - \delta_k}{\delta_k} \leq \beta$  with  $0 \leq \beta < \frac{\eta}{L}$ , we have  $\frac{\Gamma_k}{1+\beta} \leq \delta_k$ , and we obtain the following relation:

$$\mathbb{E}[\|x_{k+1} - x^*\|^2] \leq (1 - \frac{2(\eta + L)}{1 + \beta} \Gamma_k + 2L\Gamma_k + L^2\Gamma_k^2) \mathbb{E}[\|x_k - x^*\|^2] + \Gamma_k^2 \nu^2.$$

When  $\Gamma_k$  is further restricted so that  $0 < \Gamma_k \leq \frac{\eta - \beta L}{(1 + \beta)L^2}$ , we have

$$\mathbb{E}[\|x_{k+1} - x^*\|^2] \leq (1 - \frac{(\eta - \beta L)}{1 + \beta} \Gamma_k) \mathbb{E}[\|x_k - x^*\|^2] + \Gamma_k^2 \nu^2 \quad \text{for all } k \geq 0.$$

Using the preceding relation and following a similar analysis as in the proof of Proposition 4, we can show that the optimal choice of the sequence  $\{\Gamma_k^*\}$  is given by

$$\Gamma_0^* = \frac{\eta - \beta L}{2(1 + \beta)\nu^2} e_0, \quad (23)$$

$$\Gamma_k^* = \Gamma_{k-1}^* \left(1 - \frac{\eta - \beta L}{2(1 + \beta)} \Gamma_{k-1}^*\right) \quad \text{for all } k \geq 1, \quad (24)$$

where  $e_0$  is such that  $0 < e_0 \leq \frac{2\nu^2}{L^2}$ .

In the following lemma, we derive a relation between two recursive sequences, which is employed within our main convergence result for adaptive stepsizes  $\{\gamma_{k,i}\}$ .

**Lemma 4.** *Suppose that sequences  $\{\lambda_k\}$  and  $\{\gamma_k\}$  are given with the following recursive equations for all  $k \geq 0$ ,*

$$\lambda_{k+1} = \lambda_k(1 - \lambda_k), \quad \gamma_{k+1} = \gamma_k(1 - \bar{c}\gamma_k),$$

where  $\bar{c} > 0$  is a given constant and  $\lambda_0 = \bar{c}\gamma_0$ . Then for all  $k \geq 0$ ,  $\lambda_k = \bar{c}\gamma_k$ .

*Proof.* We use the induction on  $k$ . For  $k = 0$ , the relation holds since  $\lambda_0 = \bar{c}\gamma_0$ . Suppose that for some  $k \geq 0$  the relation holds. Then, we have

$$\gamma_{k+1} = \gamma_k(1 - \bar{c}\gamma_k) \Rightarrow \bar{c}\gamma_{k+1} = \bar{c}\gamma_k(1 - \bar{c}\gamma_k) \Rightarrow \bar{c}\gamma_{k+1} = \lambda_k(1 - \lambda_k) \Rightarrow \bar{c}\gamma_{k+1} = \lambda_{k+1}.$$

Hence, the result holds for  $k + 1$  implying that it holds for all  $k \geq 0$ .  $\square$

Using Lemma 4, we now present a relation between the lower and upper bound sequences given by  $\{\delta_k^*\}$  and  $\{\Gamma_k^*\}$ , respectively.

**Lemma 5.** *Suppose that the sequences  $\{\delta_k^*\}$  and  $\{\Gamma_k^*\}$  are given by relations (20)–(21) and (23)–(24), respectively, where  $0 < e_0 \leq \frac{2\nu^2}{L^2}$  and  $0 \leq \beta < \frac{\eta}{L}$ . Then, for all  $k \geq 0$ ,  $\Gamma_k^* = (1 + \beta)\delta_k^*$ .*

*Proof.* Suppose that  $\{\lambda_k\}$  is defined by the following recursive equation

$$\lambda_{k+1} = \lambda_k(1 - \lambda_k), \quad \text{for all } k \geq 0,$$

where  $\lambda_0 = \frac{(\eta - \beta L)^2}{4(1 + \beta)^2 \nu^2} e_0$ . To obtain the result, we apply Lemma 4 to sequences  $\{\lambda_k\}$  and  $\{\delta_k^*\}$ , and then to sequences  $\{\lambda_k\}$  and  $\{\Gamma_k^*\}$ . Specifically, Lemma 4 implies that  $\lambda_k = \frac{\eta - \beta L}{2} \delta_k^*$  for all  $k \geq 0$ . Invoking Lemma 4 for sequences  $\{\lambda_k\}$  and  $\{\Gamma_k^*\}$ , we have  $\lambda_k = \frac{\eta - \beta L}{2(1 + \beta)} \Gamma_k^*$ . From the preceding two relations, we conclude that  $\Gamma_k^* = (1 + \beta)\delta_k^*$  for all  $k \geq 0$ .  $\square$

The relations (20)–(21) and (23)–(24), respectively, are essentially adaptive rules for determining the best upper and lower bounds for stepsize sequences  $\{\gamma_{k,i}\}$ , where "best" corresponds to the minimizers of the associated error bounds. Having provided this intermediate result, our main result is stated next and shows the almost-sure convergence of the distributed adaptive SA (DASA) scheme.

**Theorem 1** (A class of distributed adaptive steplength SA rules). *Suppose that Assumptions 1 and 2 hold, and assume that the set  $X$  is bounded. Suppose that, for all  $i = 1, \dots, N$ , the stepsizes  $\{\gamma_{k,i}\}$  in algorithm (4) are given by the following recursive equations:*

$$\gamma_{0,i} = r_i c \frac{D^2}{\left(1 + \frac{\eta - 2c}{L}\right)^2 \nu^2}, \quad (25)$$

$$\gamma_{k,i} = \gamma_{k-1,i} \left(1 - \frac{c}{r_i} \gamma_{k-1,i}\right) \quad \text{for all } k \geq 1. \quad (26)$$

where  $D \triangleq \max_{x \in X} \|x - x_0\|$ ,  $c$  is a scalar satisfying  $c \in (0, \frac{\eta}{2})$ ,  $r_i$  is a parameter such that  $r_i \in [1, 1 + \frac{\eta - 2c}{L}]$ ,  $\eta$  is the strong monotonicity parameter of the mapping  $F$ ,  $L$  is the Lipschitz constant of  $F$ , and  $\nu$  is the upper bound defined in Assumption 2. We assume that the constant  $\nu$  is chosen large enough such that  $\nu \geq \frac{DL}{\sqrt{2}}$ . Then, the following hold:

- (a) For any  $i, j = 1, \dots, N$  and  $k \geq 0$ ,  $\frac{\gamma_{k,i}}{r_i} = \frac{\gamma_{k,j}}{r_j}$ .
- (b) Assumption 3b holds with  $\beta = \frac{\eta - 2c}{L}$ ,  $\delta_k = \delta_k^*$ ,  $\Gamma_k = \Gamma_k^*$ , and  $e_0 = D^2$ , where  $\delta_k^*$  and  $\Gamma_k^*$  are given by (20)–(21) and (23)–(24), respectively.
- (c) The sequence  $\{x_k\}$  generated by algorithm (4) converges almost surely to the unique solution of  $VI(X, F)$ .
- (d) The results of Proposition 4 hold for  $\delta_k^*$  when  $e_0 = D^2$  and  $\beta = \frac{\eta - 2c}{L}$ .

*Proof.* (a) Consider the sequence  $\{\lambda_k\}$  given by

$$\lambda_0 = \frac{c^2}{(1 + \frac{\eta - 2c}{L})^2 \nu^2} D^2, \\ \lambda_{k+1} = \lambda_k (1 - \lambda_k) \quad \text{for all } k \geq 1.$$

Since for any  $i = 1, \dots, N$ , we have  $\lambda_0 = (c/r_i) \gamma_{0,i}$ , using Lemma 4 we obtain  $\lambda_k = (c/r_i) \gamma_{k,i}$  for all  $i = 1, \dots, N$  and  $k \geq 0$ . Hence, the desired relation follows.

(b) First we show that  $\delta_k^*$  and  $\Gamma_k^*$  are well defined. Consider the relation of part (a). Let  $k \geq 0$  be arbitrarily fixed. If  $\gamma_{k,i} > \gamma_{k,j}$  for some  $i \neq j$ , then we have  $r_i > r_j$ . Therefore, the minimum possible  $\gamma_{k,i}$  is obtained with  $r_i = 1$  and the maximum possible  $\gamma_{k,i}$  is obtained with  $r_i = 1 + \frac{\eta - 2c}{L}$ . Now, consider (25)–(26). If,  $r_i = 1$ , and  $D^2$  is replaced by  $e_0$ , and  $c$  by  $\frac{\eta - \beta L}{2}$ , we get the same recursive sequence defined by (20)–(21). Therefore, since the minimum possible  $\gamma_{k,i}$  is achieved when  $r_i = 1$ , we conclude that  $\delta_k^* \leq \min_{i=1, \dots, N} \gamma_{k,i}$  for any  $k \geq 0$ . This shows that  $\delta_k^*$  is well-defined in the context of Assumption 3b. Similarly, it can be shown that  $\Gamma_k^*$  is also well-defined in the context of Assumption 3b. Now, Lemma 5 implies that  $\Gamma_k^* = (1 + \frac{\eta - 2c}{L}) \delta_k^*$  for any  $k \geq 0$ , which shows that the inequality in Assumption 3b is satisfied with  $\beta = \frac{\eta - 2c}{L}$ , where  $0 \leq \beta < \frac{\eta}{L}$  since  $0 < c \leq \frac{\eta}{2}$ .

(c) In view of Proposition 3, to show the almost-sure convergence, it suffices to show that Assumption 3 holds. Part (b) implies that Assumption 3b is satisfied by the given stepsize choices. As seen in Proposition 3 of [41], Assumption 3a holds for any positive recursive sequence  $\{\lambda_k\}$  of the form  $\lambda_{k+1} = \lambda_k (1 - a \lambda_k)$ . Since each sequence  $\gamma_{k,i}$  is a recursive sequence of this form, Assumption 3a follows from Proposition 3 in [41].

(d) It suffices to show that the hypotheses of Proposition 4 hold when  $e_0 = D^2$  and  $\beta = \frac{\eta - 2c}{L}$ . Relation  $\nu \geq \frac{DL}{\sqrt{2}}$  follows from  $\nu \geq L \sqrt{\frac{e_0}{2}}$ . Also, as mentioned in part (c), since  $0 < c \leq \frac{\eta}{2}$ , the relation  $0 \leq \beta < \frac{\eta}{L}$  holds for any choice of  $c$  within that range. Therefore, the conditions of Proposition 4 are satisfied.  $\square$

**Remark:** Theorem 1 provides a class of adaptive stepsize rules for the distributed SA algorithm (4), i.e., for any choice of parameter  $c$  such that  $0 < c \leq \frac{\eta}{2}$ , relations (25)–(26) correspond to an adaptive stepsize rule for agents  $1, \dots, N$ . Note that if  $c = \frac{\eta}{2}$ , these adaptive rules will represent the centralized adaptive scheme given by (11)–(12).  $\square$

In a distributed setting, each agent can choose its corresponding parameter  $r_i$  from the specified range  $[1, 1 + \frac{\eta - 2c}{L}]$ . This requires that all agents agree on a fixed parameter  $c$  and have a common estimate of

parameters  $\eta$  and  $L$ . Yet, this scheme does not allow complete flexibility for the agents and requires some global specification of parameters such as  $\eta$ ,  $L$ , and  $c$ . In the next section, we address the setting where the Lipschitz constant is unavailable in a global setting or when the mapping  $F$  may not be Lipschitzian are addressed.

## 4 Non-Lipschitzian mappings and local randomization

A key shortcoming of the proposed DASA scheme, given by (25)-(26), is the requirement of the Lipschitzian property of the mapping  $F$  with a known parameter  $L$ . However, in a range of problem settings, the following may arise:

- *Unavailability of a Lipschitz constant:* In many settings, either the mapping may be non-Lipschitzian or the estimation of such a constant may be problematic. It may also be that this constant may not be available across the entire population of agents.
- *Nonsmoothness in payoffs:* Suppose the Cartesian stochastic variational inequality problem represents the optimality conditions of a stochastic convex program with nonsmooth (random) objectives or the equilibrium conditions of a stochastic Nash game in which the payoff functions are expectation-valued with random nonsmooth integrands. In either setting, the integrands associated with each component's expectation are multi-valued. In such a setting, a randomization or smoothing technique applied to each agent's payoff which leads to an approximate mapping that can be shown to be Lipschitz and single-valued. The associated Lipschitz constant can be specified in terms of problem parameters and smoothing specifications, allowing us to develop a locally randomized SA algorithm for stochastic variational inequalities without Lipschitzian mappings.

In Section 4.1, we present the rudiments of our randomization approach and discuss its generalizations in Section 4.2. Finally, in Section 4.3, we present a distributed locally randomized SA scheme and provide suitable convergence theory.

### 4.1 A randomized smoothing technique

In this part, we revisit a smoothing technique that has its roots in work by Steklov [36, 37] in 1907. Over the years, it has been used by Bertsekas [3], Norkin [25] and more recently Lakshmanan and De Farias [22]. The following proposition in [3] presents this smoothing technique for a nondifferentiable convex function.

**Proposition 5.** *Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be a convex function and consider the function  $f^\epsilon(x)$*

$$f^\epsilon(x) \triangleq \mathbb{E}[f(x - \omega)],$$

*where  $\omega$  belongs to the probability space  $(\mathbb{R}^n, B_n, P)$ ,  $B_n$  is the  $\sigma$ -algebra of Borel sets of  $\mathbb{R}^n$  and  $P$  is a probability measure on  $B_n$  which is absolutely continuous with respect to Lebesgue measure restricted on  $B_n$ . Then, if  $\mathbb{E}[f(x - \omega)] < \infty$  for all  $x \in \mathbb{R}^n$ , the function  $f^\epsilon$  is everywhere differentiable.*

This technique has been employed in a number of papers such as [13, 22, 41] to transform  $f$  into a smooth function. In [22], authors consider a Gaussian distribution for the smoothing distribution and show that when function  $f$  has bounded subgradients, the smooth function  $f^\epsilon$  has Lipschitz gradients with a prescribed Lipschitz constant. A challenge in that approach is that in some situations, function  $f$  may have a restricted domain and not be defined for some realizations of the Gaussian random variable.

Motivated by this challenge, in [41], we consider the randomized smoothing technique using uniform random variables defined on an  $n$ -dimensional ball centered on origin with radius  $\epsilon > 0$ . This approach is called "locally randomized smoothing technique" and is used to establish a local smoothing SA algorithm for solving stochastic convex optimization problems in [41]. We intend to extend this smoothing technique to the regime of solving stochastic Cartesian variational inequality problems and exploit the Lipschitzian property of the approximated mapping. In the following example, we demonstrate how the smoothing technique works for a piecewise linear function.

**Example 3** (Smoothing of a convex function). *Consider the following piecewise linear function*

$$f(x) = \begin{cases} -2x - 3 & \text{for } x < -2, \\ -0.3x + 0.4 & \text{for } -2 \leq x < 3 \\ x - 3.5 & \text{for } x \geq 3. \end{cases}$$



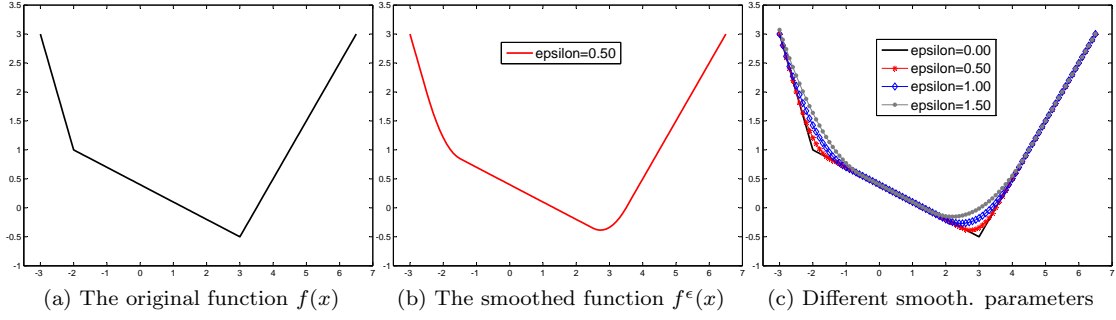


Figure 1: The smoothing technique

Suppose that  $z$  is a uniform random variable defined on  $[-\epsilon, \epsilon]$  where  $\epsilon > 0$  is a given parameter. Consider the approximation function  $f^\epsilon = \mathbb{E}[f(x + z)]$ . Proposition 5 implies that  $f^\epsilon$  is a smooth function. When  $\epsilon$  is a fixed constant satisfying  $0 < \epsilon \leq 2.5$ , the smoothed function  $f^\epsilon$  has the following form:

$$f^\epsilon(x) = \begin{cases} -2x - 3 & \text{for } x < -2 - \epsilon, \\ \frac{1}{40\epsilon} (17x^2 + 68x - 46x\epsilon + 68 - 52\epsilon + 17\epsilon^2) & \text{for } -2 - \epsilon \leq x < -2 + \epsilon, \\ -0.3x + 0.4 & \text{for } -2 + \epsilon \leq x < 3 - \epsilon, \\ \frac{1}{40\epsilon} (13x^2 - 78x + 14x\epsilon + 117 - 62\epsilon + 13\epsilon^2) & \text{for } 3 - \epsilon \leq x < 3 + \epsilon, \\ x - 3.5 & \text{for } x \geq 3 + \epsilon. \end{cases}$$

Figure 1 shows such a smoothing scheme. In Figure 1a, we observe that function  $f$  is nonsmooth at  $x = -2$  and  $x = 3$ . Figure 1b shows the approximation  $f^\epsilon$  when  $\epsilon = 0.5$ . An immediate observation is that function  $f^\epsilon$  is smooth everywhere. Furthermore, the smoothing technique perturbs  $x$  locally at all points, including points of nonsmoothness. Finally, Figure 1c shows the smoothing scheme for different values of  $\epsilon$  and illustrates the exactness of the approximation as  $\epsilon \rightarrow 0$ .

## 4.2 Locally randomized techniques

Motivated by the smoothing technique described in previous part, we introduce two distributed smoothing schemes where we simultaneously perturb the value of vectors  $x_i$  with a random vector  $z_i$  for  $i = 1, \dots, N$ . The first scheme is called a *multi-spherical randomized (MSR) scheme*, where each random vector  $z_i \in \mathbb{R}^{n_i}$  is uniformly distributed on the  $n_i$ -dimensional ball centered at the origin with radius  $\epsilon_i$ . In the second scheme, called a *multi-cubic randomized (MCR) scheme*, we let  $z_i \in \mathbb{R}^{n_i}$  be uniformly distributed on the  $n_i$ -dimensional cube centered at the origin with an edge length of  $2\epsilon_i$ .

Now, consider a mapping  $F$  that is not necessarily Lipschitz. We begin by defining an approximation  $F^\epsilon : X \rightarrow \mathbb{R}^n$  as the expectation of  $F(x)$  when  $x$  is perturbed by a random vector  $z = (z_1; \dots; z_N)$ . Specifically,  $F^\epsilon$  is given by

$$F^\epsilon(x) \triangleq \begin{pmatrix} \mathbb{E}[F_1(x + z)] \\ \vdots \\ \mathbb{E}[F_N(x + z)] \end{pmatrix} \quad \text{for all } x \in X, \quad (27)$$

where  $F_1, \dots, F_N$  are coordinate-maps of  $F$ ,  $z = (z_1; \dots; z_N)$  and the random vectors  $z_i$  are given by MSR or MCR scheme.

### 4.2.1 Multi-spherical randomized smoothing

Let us define  $B_n(x, \rho) \subset \mathbb{R}^n$  as a ball centered at a point  $x$  with a radius  $\rho > 0$ . More precisely,

$$B_n(x, \rho) \triangleq \{y \in \mathbb{R}^n \mid \|y - x\| \leq \rho\}.$$

In this scheme, assume that for all  $i = 1, \dots, N$  random vector  $z_i \in B_{n_i}(0, \epsilon_i)$  is uniformly distributed and independent with respect to random vectors  $z_j$  for  $j \neq i$ . For the approximation mapping  $F^\epsilon$  to be well-defined,  $F$  needs to be defined over the set  $X_s^\epsilon$  given by

$$X_s^\epsilon \triangleq X + \prod_{i=1}^N B_{n_i}(0, \epsilon_i).$$

This means that  $X_s^\epsilon = \{(x_1 + z_1, \dots, x_N + z_N) | x \in X, z_i \in \mathbb{R}^{n_i}, \|z_i\| \leq \epsilon_i \text{ for all } i = 1, \dots, N\}$ , where the constants  $\epsilon_i > 0$  are given values and  $\epsilon \triangleq (\epsilon_1, \dots, \epsilon_N)$ . Note that the subscript  $s$  stands for the MSR scheme. This scheme is developed based on the following assumption.

**Assumption 4.** *The mapping  $F : X_s^\epsilon \rightarrow \mathbb{R}^n$  is bounded over the set  $X_s^\epsilon$ . In particular, for every  $i = 1, \dots, N$ , there exists a constant  $C_i > 0$  such that  $\|F_i(x)\| \leq C_i$  for all  $x \in X_s^\epsilon$ .*

Under this assumption, we will show that the smoothed mapping  $F^\epsilon$  produced by the MSR scheme is Lipschitz continuous over  $X$  and we will compute its Lipschitz constant. To do so, we make use of the following lemma.

**Lemma 6.** *Let  $z \in \mathbb{R}^n$  be a random vector generated from a uniform density with zero mean over an  $n$ -dimensional ball centered at the origin with a radius  $\epsilon$ . Then, the following relation holds:*

$$\int_{\mathbb{R}^n} |p_u(z - x) - p_u(z - y)| dz \leq \kappa \frac{n!!}{(n-1)!!} \frac{\|x - y\|}{\epsilon} \quad \text{for all } x, y \in \mathbb{R}^n,$$

where  $\kappa = 1$  if  $n$  is odd and  $\kappa = \frac{2}{\pi}$  if  $n$  is even,  $n!!$  denotes double factorial of  $n$ , and  $p_u$  is the probability density function of random vector  $z$  given by

$$p_u(z) = \begin{cases} \frac{1}{c_n \epsilon^n} & \text{for } z \in B_n(0, \epsilon), \\ 0 & \text{otherwise,} \end{cases} \quad (28)$$

where  $c_n = \frac{\pi^{\frac{n}{2}}}{\Gamma(\frac{n}{2} + 1)}$ , and  $\Gamma$  is the gamma function given by

$$\Gamma\left(\frac{n}{2} + 1\right) = \begin{cases} \left(\frac{n}{2}\right)! & \text{if } n \text{ is even,} \\ \sqrt{\pi} \frac{n!!}{2^{(n+1)/2}} & \text{if } n \text{ is odd.} \end{cases}$$

*Proof.* The result is shown within the proof of Lemma 8 in the extended version of [41].  $\square$

We next provide the main result of this subsection, which establishes the Lipschitz continuity and boundedness properties of the approximation mapping  $F^\epsilon$ . It also provides the Lipschitz constant of  $F^\epsilon$  for the MSR scheme in terms of problem parameters.

**Proposition 6** (Lipschitz continuity and boundedness of  $F^\epsilon$  under the MSR scheme). *Let Assumption 4 hold and define vector  $C \triangleq (C_1, \dots, C_N)$ . Then, for any  $\epsilon = (\epsilon_1, \dots, \epsilon_N) > 0$  we have the following:*

- (a)  $F^\epsilon$  is bounded over the set  $X$ , i.e.,  $\|F^\epsilon(x)\| \leq \|C\|$  for all  $x \in X$ .
- (b)  $F^\epsilon$  is Lipschitz continuous over the set  $X$ . More precisely, we have

$$\|F^\epsilon(x) - F^\epsilon(y)\| \leq \sqrt{N} \|C\| \max_{j \in \{1, \dots, N\}} \left\{ \kappa_j \frac{n_j!!}{(n_j-1)!!} \frac{1}{\epsilon_j} \right\} \|x - y\| \quad \text{for all } x, y \in X, \quad (29)$$

where  $\kappa_j = 1$  when  $n_j$  is odd and  $\kappa_j = \frac{2}{\pi}$  when  $n_j$  is even.

*Proof.* (a) We can bound the norm of  $F^\epsilon$  as follows:

$$\|F^\epsilon(x)\| = \sqrt{\sum_{i=1}^N \|\mathbb{E}[F_i(x+z)]\|^2} \leq \sqrt{\sum_{i=1}^N \mathbb{E}[\|F_i(x+z)\|^2]} \leq \|C\|,$$

where the first inequality follows from Jensen's inequality and the second inequality is due to the boundedness property imposed on  $F$  by Assumption 4.

(b) From the definition of  $F^\epsilon$  in relation (27) we have

$$\|F^\epsilon(x) - F^\epsilon(y)\|^2 = \sum_{j=1}^N \|\mathbb{E}[F_j(x+z) - F_j(y+z)]\|^2 = \sum_{j=1}^N \|\mathbb{E}[F_j(x+z) - F_j(y+z)]\|^2.$$

We will add and subtract, sequentially, the values  $F(u)$  at the vectors  $u$  of the form  $(y_1 + z_1, \dots, y_{i-1} + z_{i-1}, x_i + z_i, \dots, x_N + z_N)$  for  $i = 2, \dots, N$ . To keep the resulting expressions in a compact form, we use the following notation. For an index set  $J \subseteq \{1, \dots, N\}$ , we let  $x_J \triangleq (x_i)_{i \in J}$  and  $x_{-J} \triangleq (x_i)_{i \in \{1, \dots, N\} - J}$ . By adding and subtracting the terms  $F_j((y+z)_{\{1, \dots, i\}}, (x+z)_{-\{1, \dots, i\}})$  for all  $i$ , from the preceding relation we obtain

$$\begin{aligned} \|F^\epsilon(x) - F^\epsilon(y)\|^2 &= \sum_{j=1}^N \left\| \underbrace{\mathbb{E}[F_j(x+z) - F_j((y+z)_{\{1\}}, (x+z)_{-\{1\}})]}_{v_1} \right. \\ &+ \underbrace{\mathbb{E}[F_j((y+z)_{\{1\}}, (x+z)_{-\{1\}}) - F_j((y+z)_{\{1,2\}}, (x+z)_{-\{1,2\}})]}_{v_2} \\ &\vdots \\ &+ \underbrace{\mathbb{E}[F_j((y+z)_{\{1, \dots, i-1\}}, (x+z)_{-\{1, \dots, i-1\}}) - F_j((y+z)_{\{1, \dots, i\}}, (x+z)_{-\{1, \dots, i\}})]}_{v_i} \\ &\vdots \\ &+ \underbrace{\mathbb{E}[F_j((y+z)_{\{1, \dots, N-2\}}, (x+z)_{-\{1, \dots, N-2\}}) - F_j((y+z)_{\{1, \dots, N-1\}}, (x+z)_{-\{1, \dots, N-1\}})]}_{v_{N-1}} \\ &\left. + \underbrace{\mathbb{E}[F_j((y+z)_{\{1, \dots, N-1\}}, (x+z)_{-\{1, \dots, N-1\}}) - F_j(y+z)]}_{v_N} \right\|^2. \end{aligned}$$

Considering the definition of the vectors  $v_1, \dots, v_N$  in the preceding relation, we have

$$\|F^\epsilon(x) - F^\epsilon(y)\|^2 = \sum_{j=1}^N \left\| \sum_{i=1}^N v_i \right\|^2 \leq N \sum_{j=1}^N \sum_{i=1}^N \|v_i\|^2,$$

where the inequality follows by the convexity of the squared-norm. By using the definitions of  $v_i$  and exchanging the order of summations in the preceding relation, we obtain

$$\begin{aligned} \|F^\epsilon(x) - F^\epsilon(y)\|^2 &\leq N \sum_{j=1}^N \underbrace{\left\| \mathbb{E}[F_j((x+z)_{\{1\}}, (x+z)_{-\{1\}}) - F_j((y+z)_{\{1\}}, (x+z)_{-\{1\}})] \right\|^2}_{\text{Term 1}} \\ &+ N \sum_{i=2}^N \sum_{j=1}^N \underbrace{\left\| \mathbb{E}[F_j((y+z)_{\{1, \dots, i-1\}}, (x+z)_{-\{1, \dots, i-1\}}) - F_j((y+z)_{\{1, \dots, i\}}, (x+z)_{-\{1, \dots, i\}})] \right\|^2}_{\text{Term } i}. \end{aligned} \quad (30)$$

Next, we derive an estimate for Term 1. From our notation, it follows that for a vector  $x$ ,  $x_{\{1\}} = x_1$ . In the interest of brevity, in the following, for a vector  $x$ , we use  $x_{-1} \triangleq x_{-\{1\}}$ . Recalling the definition of  $p_u$  in (28), we write

$$\begin{aligned} \text{Term 1} &= \sum_{j=1}^N \left\| \int_{\mathbb{R}^{n_1}} F_j(x_1 + z_1, x_{-1} + z_{-1}) p_u(z_1) dz_1 - \int_{\mathbb{R}^{n_1}} F_j(y_1 + z_1, x_{-1} + z_{-1}) p_u(z_1) dz_1 \right\|^2 \\ &= \sum_{j=1}^N \left\| \int_{\mathbb{R}^{n_1}} F_j(s_1, x_{-1} + z_{-1}) p_u(s_1 - x_1) ds_1 - \int_{\mathbb{R}^{n_1}} F_j(t_1, x_{-1} + z_{-1}) p_u(t_1 - y_1) dt_1 \right\|^2 \\ &= \sum_{j=1}^N \left\| \int_{\mathbb{R}^{n_1}} \mathbb{E}[F_j(t_1, x_{-1} + z_{-1})] (p_u(t_1 - x_1) - p_u(t_1 - y_1)) dt_1 \right\|^2, \end{aligned}$$

where in the second equality  $s_1$  and  $t_1$  are given by  $s_1 = x_1 + z_1$  and  $t_1 = y_1 + z_1$ . Using the triangle inequality and Jensen's inequality, we obtain

$$\text{Term 1} \leq \sum_{j=1}^N \left( \int_{\mathbb{R}^{n_1}} \mathbb{E}[|F_j(t_1, x_{-1} + z_{-1})|] |p_u(t_1 - x_1) - p_u(t_1 - y_1)| dt_1 \right)^2.$$

By the definition of  $F_j$  and Assumption 4, the preceding relation yields

$$\begin{aligned} \text{Term 1} &\leq \sum_{j=1}^N \left( \int_{\mathbb{R}^{n_1}} C_j |p_u(t_1 - x_1) - p_u(t_1 - y_1)| dt_1 \right)^2 \\ &\leq \left( \sum_{j=1}^N C_j^2 \right) \left( \kappa_1 \frac{n_1!!}{(n_1 - 1)!!} \frac{1}{\epsilon_1} \|x_1 - y_1\| \right)^2, \end{aligned}$$

where the last inequality is obtained using Lemma 6. Similarly, we may find estimates for the other terms in relation (30). Therefore, from relation (30) we may conclude that

$$\begin{aligned} \|F^\epsilon(x) - F^\epsilon(y)\|^2 &\leq N \left( \sum_{j=1}^N C_j^2 \right) \sum_{i=1}^N \left( \kappa_i \frac{n_i!!}{(n_i - 1)!!} \frac{1}{\epsilon_i} \|x_i - y_i\| \right)^2 \\ &\leq N \left( \sum_{j=1}^N C_j^2 \right) \left( \max_{t=1, \dots, N} \kappa_t \frac{n_t!!}{(n_t - 1)!!} \frac{1}{\epsilon_t} \right)^2 \sum_{i=1}^N \|x_i - y_i\|^2 \\ &= N \|C\|^2 \left( \max_{t=1, \dots, N} \kappa_t \frac{n_t!!}{(n_t - 1)!!} \frac{1}{\epsilon_t} \right)^2 \|x - y\|^2. \end{aligned}$$

Therefore, we have

$$\|F^\epsilon(x) - F^\epsilon(y)\| \leq \sqrt{N} \|C\| \max_{t \in \{1, \dots, N\}} \left\{ \kappa_t \frac{n_t!!}{(n_t - 1)!!} \frac{1}{\epsilon_t} \right\} \|x - y\|.$$

□

**Remark:** The MSR scheme is a generalization of the local randomization smoothing scheme presented in [41]. Note that when  $N = 1$ , the Lipschitz constant given in Proposition 6b is precisely the constant given by Lemma 8 in [41]. □

#### 4.2.2 Multi-cubic randomized smoothing scheme

We begin by defining  $C_n(x, \rho) \subset \mathbb{R}^n$  as a cube centered at a point  $x$  with the edge length  $2\rho > 0$  where the edges are along the coordinate axes. More precisely,

$$C_n(x, \rho) \triangleq \{y \in \mathbb{R}^n \mid \|y - x\|_\infty \leq \rho\}.$$

In the MCR scheme, we assume that for any  $i = 1, \dots, N$ , the random vector  $z_i$  is uniformly distributed on the set  $C_{n_i}(0, \epsilon_i)$  and is independent of the other random vectors  $z_j$  for  $j \neq i$ . For the mapping  $F$  we will assume that it is well-defined over the set  $X_c^\epsilon$  given by

$$X_c^\epsilon \triangleq X + \prod_{i=1}^N C_{n_i}(0, \epsilon_i),$$

where  $\epsilon_i > 0$  are given values and  $\epsilon \triangleq (\epsilon_1, \dots, \epsilon_N)$ , while the subscript  $c$  stands for the MCR scheme. We investigate the properties of  $F^\epsilon$  for this smoothing scheme under the following basic assumption.

**Assumption 5.** *The mapping  $F : X_c^\epsilon \rightarrow \mathbb{R}^n$  is bounded over the set  $X_c^\epsilon$ . Specifically, for every  $i = 1, \dots, N$ , there exists a constant  $C'_i > 0$  such that  $\|F_i(x)\| \leq C'_i$  for all  $x \in X_c^\epsilon$ .*

The following lemma provides a simple relation that will be important in establishing the main property of the density function used in the MCR scheme.

**Lemma 7.** *Let the vector  $p \in \mathbb{R}^m$  be such that  $0 \leq p_i \leq 1$  for all  $i = 1, \dots, m$ . Then, we have*

$$1 - \prod_{i=1}^m (1 - p_i) \leq \|p\|_1.$$

*Proof.* We use induction on  $m$  to prove this result. For  $m = 1$ , we have  $1 - \prod_{i=1}^m (1 - p_i) = p_1 = \|p\|_1$ , implying that the result holds for  $m = 1$ . Let us assume that  $1 - \prod_{i=1}^m (1 - p_i) \leq \|p\|_1$  holds for  $m$ . Therefore, we have

$$\prod_{i=1}^m (1 - p_i) \geq 1 - \sum_{i=1}^m p_i.$$

Multiplying both sides of the preceding relation by  $(1 - p_{m+1})$ , we obtain

$$\prod_{i=1}^{m+1} (1 - p_i) \geq (1 - \sum_{i=1}^m p_i)(1 - p_{m+1}) = 1 - \sum_{i=1}^{m+1} p_i + p_{m+1} \sum_{i=1}^m p_i \geq 1 - \sum_{i=1}^{m+1} p_i.$$

Hence,  $\prod_{i=1}^{m+1} (1 - p_i) \geq 1 - \sum_{i=1}^{m+1} p_i$  which implies that the result holds for  $m + 1$ . Therefore, we conclude that the result holds for any integer  $m \geq 1$ .  $\square$

The following result is crucial for establishing the properties of the approximation  $F^\epsilon$  obtained by the MCR smoothing scheme.

**Lemma 8.** *Let  $z \in \mathbb{R}^n$  be a random vector with a zero-mean uniform density over an  $n$ -dimensional cube  $\prod_{i=1}^N C_{n_i}(0, \epsilon_i)$  for  $\epsilon_i > 0$  for all  $i$ . Let the function  $p_c : \mathbb{R}^n \rightarrow \mathbb{R}$  be the probability density function of the random vector  $z$ :*

$$p_c(z) = \begin{cases} \frac{1}{2^n \prod_{i=1}^N \epsilon_i^{n_i}} & \text{for } z \in \prod_{i=1}^N C_{n_i}(0, \epsilon_i), \\ 0 & \text{otherwise.} \end{cases}$$

*Then, the following relation holds:*

$$\int_{\mathbb{R}^n} |p_c(u - x) - p_c(u - y)| du \leq \frac{\sqrt{n}}{\min_{1 \leq i \leq N} \{\epsilon_i\}} \|x - y\| \quad \text{for all } x, y \in \mathbb{R}^n.$$

*Proof.* Let  $x, y \in \mathbb{R}^n$  be arbitrary. To simplify the notation, we define sets  $S_x = \prod_{i=1}^N C_{n_i}(x_i, \epsilon_i)$  and  $S_y = \prod_{i=1}^N C_{n_i}(y_i, \epsilon_i)$ . We consider, separately, the case when the cubes  $S_x$  and  $S_y$  do not intersect, and the case when they do intersect. Before we proceed, we prove the following relation

$$S_x \cap S_y \neq \emptyset \quad \text{if and only if} \quad \|x_i - y_i\|_\infty \leq 2\epsilon_i \text{ for all } i = 1, \dots, N. \quad (31)$$

To prove relation (31), suppose that the two cubes have nonempty intersection and let  $u$  be in the intersection, i.e.,  $u \in S_x \cap S_y$ . Then, by the triangle inequality, we have for all  $i = 1, \dots, N$ ,

$$\|x_i - y_i\|_\infty \leq \|x_i - u_i\|_\infty + \|u_i - y_i\|_\infty \leq 2\epsilon_i,$$

where the last inequality follows from the fact that  $u$  belongs to each of the two cubes. Thus, when  $S_x \cap S_y \neq \emptyset$ , we have  $\|x_i - y_i\|_\infty \leq 2\epsilon_i$  for all  $i$ . Conversely, suppose now that  $\|x_i - y_i\|_\infty \leq 2\epsilon_i$  holds for all  $i = 1, \dots, N$ . Let  $\bar{u} = (x + y)/2$ , and note that by the convexity of the norm  $\|\cdot\|_\infty$ , we have

$$\|\bar{u}_i - x_i\|_\infty = \left\| \frac{y - x}{2} \right\|_\infty \leq \frac{1}{2} \|y_i - x_i\|_\infty \leq \epsilon_i \quad \text{for all } i.$$

Thus, it follows that  $\bar{u} \in S_x$ . Similarly, we find that  $\|\bar{u}_i - y_i\|_\infty \leq \epsilon_i$  for all  $i$ , which implies that  $\bar{u} \in S_y$ . Hence,  $\bar{u} \in S_x \cap S_y$ , thus showing that the two cubes have a nonempty intersection.

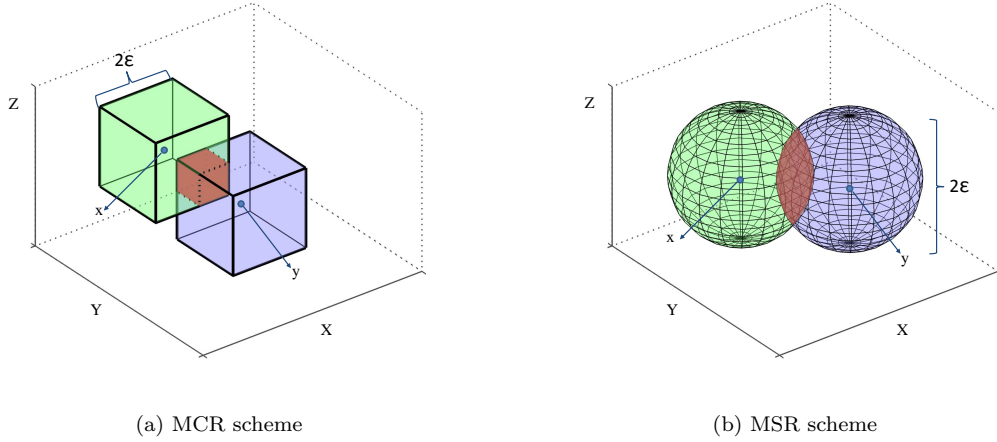


Figure 2: Calculating the Lipschitz constant in the locally randomized schemes.

We now consider the integral  $\int_{\mathbb{R}^n} |p_c(u - x) - p_c(u - y)| du$  for the cases when the cubes do not intersect and when they do intersect.

*Case 1:*  $S_x \cap S_y = \emptyset$ . In this case, we have

$$\begin{aligned} \int_{S_x} |p_c(u - x) - p_c(u - y)| du &= \int_{S_x} p_c(u - x) du, \\ \int_{S_y} |p_c(u - x) - p_c(u - y)| du &= \int_{S_y} p_c(u - y) du. \end{aligned}$$

Consequently

$$\int_{\mathbb{R}^n} |p_c(u - x) - p_c(u - y)| du = \int_{S_x} p_c(u - x) du + \int_{S_y} p_c(u - y) du = 2. \quad (32)$$

By relation (31), there must exist some index  $i^* \in \{1, \dots, N\}$  such that  $\|x_{i^*} - y_{i^*}\|_\infty > 2\epsilon_{i^*}$ . Since  $\|x - y\|_\infty \geq \|x_{i^*} - y_{i^*}\|_\infty$ , it follows that  $\frac{\|x - y\|_\infty}{\min_{1 \leq i \leq N} \{\epsilon_i\}} > 2$ . Using the relationship  $\|u\|_\infty \leq \|u\|$  between the infinity-norm and the Euclidean norm, we obtain  $\frac{\|x - y\|}{\min_{1 \leq i \leq N} \{\epsilon_i\}} > 2$ . Therefore, using (32), we have

$$\int_{\mathbb{R}^n} |p_c(u - x) - p_c(u - y)| du < \frac{1}{\min_{1 \leq i \leq N} \{\epsilon_i\}} \|x - y\|. \quad (33)$$

Case 2:  $S_x \cap S_y \neq \emptyset$ . Then, we may decompose the integral as follows:

$$\begin{aligned} \int_{\mathbb{R}^n} |p_c(u-x) - p_c(u-y)| du &= \int_{S_x \cap S_y} |p_c(u-x) - p_c(u-y)| du + \int_{S_x^c \cap S_y^c} |p_c(u-x) - p_c(u-y)| du \\ &\quad + \int_{S_x \setminus S_y} |p_c(u-x) - p_c(u-y)| du + \int_{S_y \setminus S_x} |p_c(u-x) - p_c(u-y)| du. \end{aligned}$$

Note that the first two integrals on the right hand side of the preceding equality are zero since  $p_c(u-x) = p_c(u-y)$  in the corresponding regions. Figure 2a illustrates this observation<sup>2</sup>. Therefore, we have

$$\int_{\mathbb{R}^n} |p_c(u-x) - p_c(u-y)| du = \int_{S_x \setminus S_y} p_c(u-x) du + \int_{S_y \setminus S_x} p_c(u-y) du = 2 \frac{1}{2^n \prod_{i=1}^N \epsilon_i^{n_i}} \int_{S_x \setminus S_y} du.$$

Note that the value  $2^n \prod_{i=1}^N \epsilon_i^{n_i}$  is the volume of the cube  $S_x$ , denoted by  $\text{vol}(S_x)$ . Similarly, the integral  $\int_{S_x \setminus S_y} du$  is equal to the volume of the set  $S_x \setminus S_y$ . Thus, we can write

$$\int_{\mathbb{R}^n} |p_c(u-x) - p_c(u-y)| du = 2 \frac{\text{vol}(S_x \setminus S_y)}{\text{vol}(S_x)} = 2 \frac{\text{vol}(S_x) - \text{vol}(S_x \cap S_y)}{\text{vol}(S_x)} = 2 \left( 1 - \frac{\text{vol}(S_x \cap S_y)}{\text{vol}(S_x)} \right).$$

It can be seen that

$$\text{vol}(S_x \cap S_y) = \prod_{i=1}^N \prod_{j=1}^{n_i} (2\epsilon_i - |x_i(j) - y_i(j)|),$$

where  $w(j)$  denotes the  $j$ -th coordinate value of a vector  $w$ . Therefore, from the preceding two relations and  $\text{vol}(S_x) = 2^n \prod_{i=1}^N \epsilon_i^{n_i}$  we find that

$$\begin{aligned} \int_{\mathbb{R}^n} |p_c(u-x) - p_c(u-y)| du &= 2 \left( 1 - \frac{1}{2^n \prod_{i=1}^N \epsilon_i^{n_i}} \left( \prod_{i=1}^N \prod_{j=1}^{n_i} (2\epsilon_i - |x_i(j) - y_i(j)|) \right) \right) \\ &= 2 \left( 1 - \prod_{i=1}^N \prod_{j=1}^{n_i} \left( 1 - \frac{|x_i(j) - y_i(j)|}{2\epsilon_i} \right) \right). \end{aligned} \quad (34)$$

Since the cubes  $S_x$  and  $S_y$  do intersect, by relation (31) there must hold  $\|x_i - y_i\|_\infty \leq 2\epsilon_i$  for all  $i$ . Hence,  $0 \leq \frac{|x_i(j) - y_i(j)|}{2\epsilon_i} \leq 1$  for all  $i$ . Now, invoking Lemma 7, from (34) we obtain

$$\int_{\mathbb{R}^n} |p_c(u-x) - p_c(u-y)| du \leq 2 \sum_{i=1}^N \sum_{j=1}^{n_i} \frac{|x_i(j) - y_i(j)|}{2\epsilon_i} = \sum_{i=1}^N \frac{\|x_i - y_i\|_1}{\epsilon_i} \leq \sum_{i=1}^N \frac{\sqrt{n_i}}{\epsilon_i} \|x_i - y_i\|,$$

where in the last inequality we used the relation between  $\|\cdot\|_1$  and the Euclidean norm. Using Hölder's inequality, we have

$$\sum_{i=1}^N \frac{\sqrt{n_i}}{\epsilon_i} \|x_i - y_i\| \leq \sqrt{\sum_{i=1}^N \frac{n_i}{\epsilon_i^2}} \|x - y\| \leq \frac{\sqrt{n}}{\min_{1 \leq i \leq N} \{\epsilon_i\}} \|x - y\|,$$

implying that

$$\int_{\mathbb{R}^n} |p_c(u-x) - p_c(u-y)| du \leq \frac{\sqrt{n}}{\min_{1 \leq i \leq N} \{\epsilon_i\}} \|x - y\|. \quad (35)$$

By combining (37), (33), and (35), and using the fact  $n \geq 1$ , we obtain the desired result.  $\square$

Analogous to Proposition 6, the next proposition derives the Lipschitz constant and boundedness properties of the approximation  $F^\epsilon$  under the MCR scheme.

<sup>2</sup>Figure 2b provides a similar graphic for the MSR scheme.



**Proposition 7** (Lipschitz continuity and boundedness of  $F^\epsilon$  under the MCR scheme). *Let Assumption 5 hold and define vector  $C' \triangleq (C'_1, \dots, C'_N)$ . Then, for any  $\epsilon = (\epsilon_1, \dots, \epsilon_N) > 0$  we have the following:*

- (a)  $F^\epsilon$  is bounded over the set  $X$ , i.e.,  $\|F^\epsilon(x)\| \leq \|C'\|$  for all  $x \in X$ .
- (b)  $F^\epsilon$  is Lipschitz over the set  $X$ . More precisely, we have

$$\|F^\epsilon(x) - F^\epsilon(y)\| \leq \frac{\sqrt{n}\|C'\|}{\min_{j=1, \dots, N} \{\epsilon_j\}} \|x - y\| \quad \text{for all } x, y \in X. \quad (36)$$

*Proof.* (a) This result can be shown in a similar fashion to the proof of Proposition 6a.

(b) Since the random vector  $z_i$  is uniformly distributed on the set  $C_{n_i}(0, \epsilon_i)$  for each  $i = 1, \dots, N$ , the random vector  $z = (z_1; \dots; z_N)$  is uniformly distributed on the set  $\prod_{i=1}^N C_{n_i}(0, \epsilon_i)$ . By the definition of the approximation  $F^\epsilon$  in (27), it follows that for any  $x, y \in X$ ,

$$\begin{aligned} \|F^\epsilon(x) - F^\epsilon(y)\| &= \left\| \int_{\mathbb{R}^n} F(x+z)p_c(z)dz - \int_{\mathbb{R}^n} F(y+z)p_c(z)dz \right\| \\ &= \left\| \int_{\mathbb{R}^n} F(u)p_c(u-x)du - \int_{\mathbb{R}^n} F(v)p_c(v-y)dv \right\| \\ &= \left\| \int_{\mathbb{R}^n} F(u)(p_c(u-x) - p_c(u-y))du \right\| \\ &\leq \int_{\mathbb{R}^n} \|F(u)\| |p_c(u-x) - p_c(u-y)| du, \end{aligned}$$

where in the second equality we let  $u = x + z$  and  $v = y + z$ , while the inequality follows from the triangle inequality. Invoking Assumption 5 we obtain

$$\|F^\epsilon(x) - F^\epsilon(y)\| \leq \|C'\| \int_{\mathbb{R}^n} |p_c(u-x) - p_c(u-y)| du. \quad (37)$$

The desired relation follows from relation (37) and Lemma 8.  $\square$

### 4.3 A distributed locally randomized SA scheme

The locally randomized schemes presented in Section 4.2 facilitate the construction of a distributed locally randomized SA scheme. Consider the Cartesian stochastic variational inequality problem  $\text{VI}(X, F^\epsilon)$  given in (27) where the mapping  $F$  is not necessarily Lipschitz. In this section, we assume that the conditions of the MSR scheme are satisfied, i.e., for all  $i = 1, \dots, N$ , the random vector  $z_i$  is uniformly distributed over the set  $\in B_{n_i}(0, \epsilon_i)$  independently from the other random vectors  $z_j$  for  $j \neq i$ , and the mapping  $F$  in (2) is defined over the set  $X_s^\epsilon$ . Let the sequence  $\{x_k\}$  be given by

$$x_{k+1,i} = \Pi_{X_i}(x_{k,i} - \gamma_{k,i} \Phi_i(x_k + z_k, \xi_k)), \quad (38)$$

for all  $k \geq 0$  and  $i = 1, \dots, N$ , where  $\gamma_{k,i} > 0$  denotes the stepsize of the  $i$ -th agent at iteration  $k$ ,  $x_k = (x_{k,1}; x_{k,2}; \dots; x_{k,N})$ , and  $z_k = (z_{k,1}; z_{k,2}; \dots; z_{k,N})$ . The following proposition proves the almost-sure convergence of the iterates generated by algorithm (38) to the solution of the approximation  $\text{VI}(X, F^\epsilon)$ . In this result, we proceed to show that the approximation does indeed satisfy the assumptions of Proposition 3 and convergence can then be immediately claimed. We define  $\mathcal{F}'_k$ , the history of the method up to time  $k$ , as

$$\mathcal{F}'_k \triangleq \{x_0, z_0, \xi_0, z_1, \xi_1, \dots, z_{k-1}, \xi_{k-1}\},$$

for  $k \geq 1$  and  $\mathcal{F}'_0 = \{x_0\}$ . We assume that, at any iteration  $k$ , the vectors  $z_k$  and  $\xi_k$  in (38) are independent given the history  $\mathcal{F}'_k$ .

**Proposition 8** (Almost-sure convergence of locally randomized DASA scheme). *Let Assumptions 1a, 3, and 4 hold, and suppose that mapping  $F$  is strongly monotone on the set  $X_s^\epsilon$  with a constant  $\eta > 0$ . Also, assume that, for each  $i = 1, \dots, N$ , there exists a constant  $\nu_i > 0$  such that*

$$\mathbb{E}[\|\Phi_i(x_k + z_k, \xi_k) - F_i(x_k + z_k)\|^2 \mid \mathcal{F}'_k] \leq \nu_i^2 \quad \text{a.s. for all } k. \quad (39)$$

Then, the sequence  $\{x_k\}$  generated by algorithm (38) converges almost surely to the unique solution of  $VI(X, F^\epsilon)$ .

*Proof.* Define random vector  $\xi' \triangleq (z_1; z_2; \dots; z_N; \xi)$ , allowing us to rewrite algorithm (38) as follows:

$$\begin{aligned} x_{k+1,i} &= \Pi_{X_i} (x_{k,i} - \gamma_{k,i} (F_i^\epsilon(x_k) + w'_{k,i})), \\ w'_{k,i} &\triangleq \Phi_i(x_k + z_k, \xi_k) - F_i^\epsilon(x_k). \end{aligned} \quad (40)$$

To prove convergence of the iterates produced by (40), it suffices to show that the conditions of Proposition 3 are satisfied for the set  $X$ , the mapping  $F^\epsilon$ , and the stochastic errors  $w'_{k,i}$ .

(i) Since Assumption 4 holds, Proposition 6b implies that the mapping  $F^\epsilon$  is Lipschitz over the set  $X$  with the constant  $\sqrt{N}\|C\| \max_{1 \leq j \leq N} \{\kappa_j \frac{n_j!!}{(n_j-1)!!} \frac{1}{\epsilon_j}\}$ . Thus, Assumption 1b holds for the mapping  $F^\epsilon$ .

(ii) Next, we show that the mapping  $F^\epsilon$  is strongly monotone over  $X$ . Since the mapping  $F$  is strongly monotone over the set  $X_s^\epsilon$  with a constant  $\eta > 0$ , for any  $u, v \in X_s^\epsilon$ , we have

$$(u - v)^T (F(u) - F(v)) \geq \eta \|u - v\|^2.$$

Therefore, for any  $x, y \in X$  and any realization of the random vector  $z$ , the vectors  $x + z$  and  $y + z$  belong to the set  $X_s^\epsilon$ . Consequently, by defining  $u \triangleq x + z$  and  $v \triangleq y + z$ , respectively, and noting that  $u - v = x - y$ , from the previous relation we obtain

$$(x - y)^T (F(x + z) - F(y + z)) \geq \eta \|x - y\|^2.$$

Taking expectations on both sides, it follows that

$$(x - y)^T (\mathbb{E}[F(x + z)] - \mathbb{E}[F(y + z)]) \geq \eta \|x - y\|^2,$$

which implies that  $F^\epsilon$  is strongly monotone over the set  $X$  with the constant  $\eta$ .

(iii) The last step of the proof entails showing that the stochastic errors  $w'_k \triangleq (w_{k,1}; w_{k,2}; \dots; w_{k,N})$  are well-defined, i.e.,  $\mathbb{E}[w'_k | \mathcal{F}'_k] = 0$  and that Assumption 2 holds with respect to the stochastic error  $w'_k$ . Consider the definition of  $w'_{k,i}$  in (40). Taking conditional expectations on both sides, we have for all  $i = 1, \dots, N$

$$\mathbb{E}[w'_{k,i} | \mathcal{F}'_k] = \mathbb{E}_{z,\xi}[\Phi_i(x_k + z_k, \xi_k)] - F_i^\epsilon(x_k) = \mathbb{E}[F_i(x_k + z_k)] - F_i^\epsilon(x_k) = F_i^\epsilon(x_k) - F_i^\epsilon(x_k) = 0,$$

where the last equality is obtained using the definition of  $F^\epsilon$  in (27). Consequently, it suffices to show that the condition of Assumption 2 holds. This may be expressed as follows:

$$\mathbb{E}[\|w'_k\|^2 | \mathcal{F}'_k] = \mathbb{E}\left[\sum_{i=1}^N \|w'_{k,i}\|^2 | \mathcal{F}'_k\right] = \mathbb{E}_{z,\xi}\left[\sum_{i=1}^N \|\Phi_i(x_k + z_k, \xi_k) - F_i^\epsilon(x_k)\|^2 | \mathcal{F}'_k\right].$$

By adding and subtracting  $F_i(x_k + z_k)$  we obtain

$$\begin{aligned} \mathbb{E}[\|w'_k\|^2 | \mathcal{F}'_k] &\leq 2\mathbb{E}_{z,\xi}\left[\sum_{i=1}^N (\|\Phi_i(x_k + z_k, \xi_k) - F_i(x_k + z_k)\|^2 + \|F_i(x_k + z_k) - F_i^\epsilon(x_k)\|^2) | \mathcal{F}'_k\right] \\ &= 2\sum_{i=1}^N \mathbb{E}[\mathbb{E}[\|\Phi_i(x_k + z_k, \xi_k) - F_i(x_k + z_k)\|^2 | \mathcal{F}'_k, z_k] | \mathcal{F}'_k] \\ &\quad + 2\sum_{i=1}^N \mathbb{E}[(\|F_i(x_k + z_k)\|^2 - \|F_i^\epsilon(x_k)\|^2) | \mathcal{F}'_k], \end{aligned}$$

where the last term is obtained from the following relation:

$$\mathbb{E}[F_i(x_k + z_k)^T F_i^\epsilon(x_k) | \mathcal{F}'_k] = \mathbb{E}[F_i(x_k + z_k)^T F_i^\epsilon(x_k) | x_k] = \|F_i^\epsilon(x_k)\|^2.$$

Using the assumption on the errors given in (39), we further obtain

$$\mathbb{E}[\|w'_k\|^2 \mid \mathcal{F}'_k] \leq 2 \sum_{i=1}^N \nu_i^2 + 2 \sum_{i=1}^N \mathbb{E}[(\|F_i(x_k + z_k)\|^2 - \|F_i^\epsilon(x_k)\|^2) \mid \mathcal{F}'_k]. \quad (41)$$

Furthermore, we have

$$\sum_{i=1}^N \mathbb{E}[(\|F_i(x_k + z_k)\|^2 - \|F_i^\epsilon(x_k)\|^2) \mid \mathcal{F}'_k] \leq \sum_{i=1}^N \mathbb{E}[\|F_i(x_k + z_k)\|^2 \mid \mathcal{F}'_k] \leq C^2, \quad (42)$$

where we use the fact  $x_k + z_k \in X_s^\epsilon$  and the assumption that  $F_i$  is uniformly bounded over the set  $X_s^\epsilon$  (cf. Assumption 4). Relations (41)–(42) imply that the stochastic errors  $\{w'_k\}$  satisfy Assumption 2. Thus, the conditions of Proposition 3 are satisfied for the set  $X$ , the mapping  $F^\epsilon$ , and the stochastic errors  $w'_{k,i}$  and the convergence result follows.  $\square$

The distributed locally randomized SA scheme produces a solution that is an approximation to the true solution. A natural question is whether the sequence of approximations tends to the solution of  $\text{VI}(X, F)$  as  $\epsilon$ , the size of the support of the randomization, tends to zero. The following proposition resolves this question in the affirmative.

**Proposition 9.** *Let Assumption 1a hold, and suppose that mapping  $F$  is a continuous and strongly monotone over the set  $X_s^\epsilon$ . Let  $x^\epsilon$  and  $x^*$  denote the solution of  $\text{VI}(X, F^\epsilon)$  and  $\text{VI}(X, F)$ , respectively. Then  $x^\epsilon \rightarrow x^*$  when  $\epsilon \rightarrow 0$ .*

*Proof.* As showed in the proof of Proposition 8,  $F^\epsilon$  is also strongly monotone over the set  $X$  with constant  $\eta$ . Since set  $X$  is assumed to be closed and convex, the definition of  $X_s^\epsilon$  implies that  $X_s^\epsilon$  is also closed and convex. Thus, the existence and uniqueness of the solution to  $\text{VI}(X, F)$ , as well as  $\text{VI}(X, F^\epsilon)$ , is guaranteed by Theorem 2.3.3 of [11].

Let  $\epsilon = (\epsilon_1, \epsilon_2, \dots, \epsilon_N)$  with  $\epsilon_i > 0$  for all  $i$  be arbitrary, and let  $x^\epsilon$  denote the solution to  $\text{VI}(X, F^\epsilon)$ . Let  $x^*$  be the solution to  $\text{VI}(X, F)$ . Thus, since  $x^\epsilon$  is the solution to  $\text{VI}(X, F^\epsilon)$ , we have  $(x^* - x^\epsilon)^T F^\epsilon(x^\epsilon) \geq 0$ . Similarly, since  $x^*$  is the solution to  $\text{VI}(X, F)$ , we have  $(x^\epsilon - x^*)^T F(x^*) \geq 0$ . Adding the preceding two inequalities, we obtain for any  $k \geq 0$ ,

$$(x^* - x^\epsilon)^T (F^\epsilon(x^\epsilon) - F(x^*)) \geq 0.$$

Adding and subtracting the term  $F^\epsilon(x^*)$ , we have

$$(x^* - x^\epsilon)^T (F^\epsilon(x^\epsilon) - F^\epsilon(x^*)) + (x^* - x^\epsilon)^T (F^\epsilon(x^*) - F(x^*)) \geq 0,$$

implying that

$$(x^* - x^\epsilon)^T (F^\epsilon(x^*) - F(x^*)) \geq (x^* - x^\epsilon)^T (F^\epsilon(x^*) - F^\epsilon(x^\epsilon)) \geq \eta \|x^* - x^\epsilon\|^2,$$

where the last inequality follows by the strong monotonicity of the mapping  $F^\epsilon$ . By invoking the Cauchy-Schwartz inequality, we obtain

$$\|F^\epsilon(x^*) - F(x^*)\| \geq \eta \|x^* - x^\epsilon\|. \quad (43)$$

Next, we show that  $\lim_{\epsilon \rightarrow 0} F^\epsilon(x^*) = F(x^*)$ . By the definition of  $F^\epsilon$  and Jensen's inequality, we have

$$\|F^\epsilon(x^*) - F(x^*)\| = \|\mathbb{E}[F(x^* + z) - F(x^*)]\| \leq \mathbb{E}[\|F(x^* + z) - F(x^*)\|]. \quad (44)$$

Then, the expectation on the right-hand side can be expressed as follows:

$$\begin{aligned} \mathbb{E}[\|F(x^* + z) - F(x^*)\|] &= \int_{\mathbb{R}^{n_1}} \dots \int_{\mathbb{R}^{n_N}} \|F(x^* + z) - F(x^*)\| \left( \prod_{i=1}^N p_u(z_i) \right) dz_1 \dots dz_N \\ &= \int_{B_{n_1}(0, \epsilon_1)} \dots \int_{B_{n_N}(0, \epsilon_N)} \|F(x^* + z) - F(x^*)\| \left( \prod_{i=1}^N p_u(z_i) \right) dz_1 \dots dz_N, \end{aligned} \quad (45)$$

where the second equality is a consequence of the definition of the random vector  $z$ . Let  $\delta > 0$  be an arbitrary fixed number. By the continuity of  $F$  over  $X_s^\epsilon$ , there exists a  $\delta' > 0$ , such that if  $\|(x^* + z) - x^*\| \leq \delta'$ , then  $\|F(x^* + z) - F(x^*)\| \leq \delta$ . Therefore, for all  $\epsilon = (\epsilon_1, \epsilon_2, \dots, \epsilon_N)$  with  $\|\epsilon\| \leq \delta'$  we have  $\|z\| \leq \|\epsilon\| \leq \delta'$  for  $z \in \prod_{i=1}^N B_{n_i}(0, \epsilon_i)$ , which is equivalent to  $\|(x^* + z) - x^*\| \leq \delta'$ . Hence,  $\|F(x^* + z) - F(x^*)\| \leq \delta$  for all  $z \in \prod_{i=1}^N B_{n_i}(0, \epsilon_i)$  with  $\epsilon_i$  such that  $\|\epsilon\| \leq \delta'$ . Thus, using (44) and (45), for any  $\epsilon = (\epsilon_1, \dots, \epsilon_N)$  with  $\|\epsilon\| \leq \delta'$ , we have

$$\|F^\epsilon(x^*) - F(x^*)\| \leq \delta \int_{B_{n_1}(0, \epsilon_{k,1})} \dots \int_{B_{n_N}(0, \epsilon_{k,N})} \left( \prod_{i=1}^N p_u(z_i) \right) dz_{k,1} \dots dz_{k,N} = \delta.$$

Since  $\delta > 0$  was arbitrary, we conclude that  $\lim_{\epsilon \rightarrow 0} \|F_k(x^*) - F(x^*)\| = 0$ . Therefore, taking limits on both sides of inequality (43), we obtain  $\lim_{\epsilon \rightarrow 0} \|x^* - x^{\epsilon_k}\| = 0$ .  $\square$

**Remark:** Note that the results of Propostion 8 and Proposition 9 hold when the random vector  $z$  fits the conditions of the MCR scheme.

## 5 Numerical results

In this section, we report the results of our numerical experiments on two sets of test problems. Of these, the first is a stochastic bandwidth-sharing problem in communication networks (Sec. 5.1), while the second is a stochastic Nash-Cournot game (Sec. 5.2). In each instance, we compare the performance of the distributed adaptive stepsize SA scheme (DASA) given by (25)–(26) with that of SA schemes with harmonic stepsize sequences (HSA), where agents use the stepsize  $\frac{\theta}{k}$  at iteration  $k$ . More precisely, we consider three different values of the parameter  $\theta$ , i.e.,  $\theta = 0.1, 1$ , and  $10$ . This diversity of choices allows us to observe the sensitivity of the HSA scheme to different settings of the parameters. In the context of Nash-Cournot games, we use the distributed locally randomized SA scheme described in Sec. 4.3 with the MSR and MCR techniques. In each instance, we conduct a sensitivity analysis where we consider 12 different parameter settings, categorized into 4 sets. In each set, one parameter is changed while other parameters are maintained as fixed. We provide 90% confidence intervals of the mean squared error for each of the 12 settings. Our experiments have been done using Matlab 7.12.

### 5.1 A bandwidth-sharing problem in computer networks

We consider a communication network where users compete for the bandwidth. Such a problem can be captured by an optimization framework (cf. [6]). Motivated by this model, we consider a network with 16 nodes, 20 links and 5 users. Figure 3 shows the configuration of this network. Users have access to different routes as shown in Figure 3. For example, user 1 can access routes 1, 2, and 3. Each user is characterized by a cost function. Additionally, there is a congestion cost function that depends on the aggregate flow. More specifically, the cost function user  $i$  with flow rate (bandwidth)  $x_i$  is defined by

$$f_i(x_i, \xi_i) \triangleq - \sum_{r \in \mathcal{R}(i)} \xi_i(r) \log(1 + x_i(r)),$$

for  $i = 1, \dots, 5$ , where  $x \triangleq (x_1; \dots; x_5)$  is the flow decision vector of the users,  $\xi \triangleq (\xi_1; \dots; \xi_5)$  is a random parameter corresponding to the different users,  $\mathcal{R}(i) = \{1, 2, \dots, n_i\}$  is the set of routes assigned to the  $i$ -th user,  $x_i(r)$  and  $\xi_i(r)$  are the  $r$ -th element of the decision vector  $x_i$  and the random vector  $\xi_i$ , respectively. We assume that  $\xi_i(r)$  is drawn from a uniform distribution for each  $i$  and  $r$ . More precisely,  $\xi_1(1)$ ,  $\xi_1(2)$ , and  $\xi_1(3)$  are i.i.d. and uniformly distributed in  $[1 - 0.1, 1 + 0.1]$ ,  $\xi_2(1)$  and  $\xi_2(2)$  are i.i.d. and uniformly distributed in  $[1.4 - 0.2, 1.4 + 0.2]$ ,  $\xi_3(1)$  and  $\xi_4(1)$  are i.i.d. and uniformly distributed in  $[0.8 - 0.05, 0.8 + 0.05]$  and  $[1.6 - 0.2, 1.6 + 0.2]$ , respectively, and  $\xi_5(1)$  and  $\xi_5(2)$  are i.i.d and uniformly distributed in  $[1.2 - 0.1, 1.2 + 0.1]$ .

The links have limited capacities, which are given by

$$b = (10; 15; 15; 20; 10; 10; 20; 30; 25; 15; 20; 15; 10; 10; 15; 15; 20; 20; 25; 40).$$

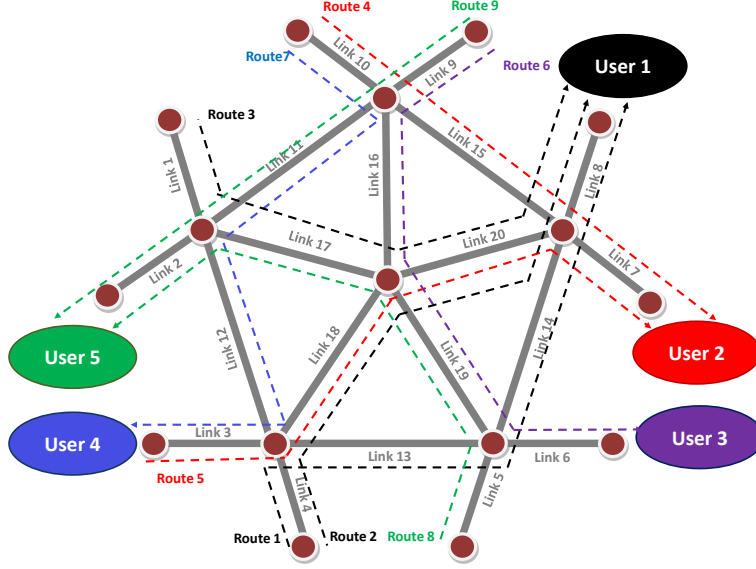


Figure 3: The bandwidth-sharing problem – the network

We may define the routing matrix  $A$  that describes the relation between set of routes  $\mathcal{R} = \{1, 2, \dots, 9\}$  and set of links  $\mathcal{L} = \{1, 2, \dots, 20\}$ . Assume that  $A_{lr} = 1$  if route  $r \in \mathcal{R}$  goes through link  $l \in \mathcal{L}$  and  $A_{lr} = 0$  otherwise. Using this matrix, the capacity constraints of the links can be described by  $Ax \leq b$ .

We formulate this model as a stochastic optimization problem given by

$$\begin{aligned} & \text{minimize} && \sum_{i=1}^N \mathbb{E}[f_i(x_i, \xi_i)] + c(x) \\ & \text{subject to} && Ax \leq b \\ & && x \geq 0, \end{aligned} \tag{46}$$

where  $c(x)$  is the network congestion cost. We consider this cost of the form  $c(x) = \|Ax\|^2$ . Problem (46) is a convex optimization problem and the optimality conditions can be stated as a variational inequality given by  $\nabla f(x^*)^T(x - x^*) \geq 0$ , where  $f(x) \triangleq \sum_{i=1}^N \mathbb{E}[f_i(x_i, \xi_i)] + c(x)$ . Using our notation in Sec. 2.2, we have

$$F(x) \triangleq \nabla f(x) = - \left( \frac{\bar{\xi}_1(1)}{1 + x_1(1)}; \dots; \frac{\bar{\xi}_i(r_i)}{1 + x_i(r_i)}; \dots; \frac{\bar{\xi}_5(2)}{1 + x_5(2)} \right) + 2A^T Ax,$$

where  $\bar{\xi}_i(r_i) \triangleq \mathbb{E}[\xi_i(r_i)]$  for any  $i = 1, \dots, 5$  and  $r_i = 1, \dots, n_i$ . We now show that the mapping  $F$  is Lipschitz and strongly monotone. Using the preceding relation, triangle inequality, and Cauchy-Schwartz inequality, for any  $x, y \in X \triangleq \{x \in \mathbb{R}^N | Ax \leq b, x \geq 0\}$ , we have

$$\begin{aligned} & \|F(x) - F(y)\| \\ &= \left\| - \left( \bar{\xi}_1(1) \left( \frac{1}{1 + x_1(1)} - \frac{1}{1 + y_1(1)} \right); \dots; \bar{\xi}_5(2) \left( \frac{1}{1 + x_5(2)} - \frac{1}{1 + y_5(2)} \right) \right) + 2A^T A(x - y) \right\| \\ &\leq \left\| \left( \bar{\xi}_1(1) \frac{x_1(1) - y_1(1)}{(1 + x_1(1))(1 + y_1(1))}; \dots; \bar{\xi}_5(2) \frac{x_5(2) - y_5(2)}{(1 + x_5(2))(1 + y_5(2))} \right) \right\| + 2\|A^T A\| \|x - y\|. \end{aligned}$$

Using nonnegativity constraints, from the preceding relation we obtain

$$\|F(x) - F(y)\| \leq \max_{i, r_i} \bar{\xi}_i(r_i) \|x - y\| + 2\|A^T A\| \|x - y\| = \left( \max_{i, r_i} \bar{\xi}_i(r_i) + 2\|A^T A\| \right) \|x - y\|,$$

implying that  $F$  is Lipschitz with constant  $\max_{i,r_i} \bar{\xi}_i(r_i) + 2\|A^T A\|$ . To show the monotonicity of  $F$ , we write

$$\begin{aligned}
& (F(x) - F(y))^T (x - y) \\
&= \left( \left( \bar{\xi}_1(1) \frac{x_1(1) - y_1(1)}{(1 + x_1(1))(1 + y_1(1))}; \dots; \bar{\xi}_5(2) \frac{x_5(2) - y_5(2)}{(1 + x_5(2))(1 + y_5(2))} \right) + 2A^T A(x - y) \right)^T (x - y) \\
&= \sum_{i,r} \bar{\xi}_i(r) \frac{(x_i(r) - y_i(r))^2}{(1 + x_i(r))(1 + y_i(r))} + 2(x - y)^T (A^T A)(x - y) \\
&\geq \frac{\min_{i,r_i} \bar{\xi}_i(r_i)}{(1 + \max_l b(l))^2} \|x - y\|^2 + 2(x - y)^T (A^T A)(x - y) \\
&= (x - y)^T \left( \frac{\min_{i,r_i} \bar{\xi}_i(r_i)}{(1 + \max_l b(l))^2} \mathbf{I}_N + 2A^T A \right) (x - y).
\end{aligned}$$

Our choice of matrix  $A$  is such that  $A^T A$  is positive definite. Thus, the preceding relation implies that  $F$  is strongly monotone with parameter

$$\eta = \frac{\min_{i,r} \bar{\xi}_i(r_i)}{(1 + \max_l b(l))^2} + 2\lambda_{\min}(A^T A),$$

where  $\lambda_{\min}(A^T A)$  is the minimum eigenvalue of the matrix  $A^T A$ .

### 5.1.1 Specification of parameters

In this experiment, the optimal solution  $x^*$  of the problem (46) is calculated by sample average approximation (SAA) method using the nonlinear programming solver `knitro` [5]. Our goal lies in comparing the performance of the DASA scheme given by (25)–(26) with that of SA schemes using harmonic stepsize sequences of the form  $\gamma_k = \frac{\theta}{k}$ , referred to as HSA schemes. We consider three values for  $\theta$  and observe the performance of HSA scheme in each case. To calculate the stepsize sequence in DASA scheme, other than  $\eta$  and  $L$  obtained in the previous part, parameters  $c$ ,  $r_i$ ,  $D$ , and  $\nu$  need to be evaluated. We assume that  $c = \frac{\eta}{4}$  and  $r_i$  is uniformly drawn from the interval  $[1, 1 + \frac{\eta - 2c}{L}]$  for each user. We let the starting point of all SA schemes be zero, i.e.,  $x_0 = 0$ . Thus,  $D = \max_{x \in X} \|x\|$ . Since the routing matrix  $A$  has binary entries, from  $Ax \leq b$ , one may conclude that  $\sqrt{N} \max_l b(l)$  can be chosen as  $D$ . To calculate  $\nu$ , for any  $k \geq 0$  we have

$$\begin{aligned}
\mathbb{E}[\|w_k\|^2 \mid \mathcal{F}_k] &= \mathbb{E}[\|\Phi(x_k, \xi_k) - F(x_k)\|^2 \mid \mathcal{F}_k] \\
&= \mathbb{E} \left[ \left\| \left( \frac{\xi_{k,1}(1) - \bar{\xi}_{k,1}(1)}{1 + x_{k,1}(1)}; \dots; \frac{\xi_{k,5}(2) - \bar{\xi}_{k,5}(2)}{1 + x_{k,5}(2)} \right) \right\|^2 \mid \mathcal{F}_k \right] \\
&= \mathbb{E} \left[ \sum_{i=1}^N \sum_{r=1}^{n_i} \left( \frac{\xi_{k,i}(r) - \bar{\xi}_{k,i}(r)}{1 + x_{k,i}(r)} \right)^2 \mid \mathcal{F}_k \right] \\
&= \sum_{i=1}^N \sum_{r=1}^{n_i} \frac{\text{var}(\xi_{k,i}(r))}{(1 + x_{k,i}(r))^2} \\
&\leq \sum_{i=1}^N \sum_{r=1}^{n_i} \text{var}(\xi_{k,i}(r)),
\end{aligned}$$

where the last inequality is obtained using  $x_{k,i}(r) \geq 0$ . Thus,  $\sqrt{\sum_{i=1}^N \sum_{r=1}^{n_i} \text{var}(\xi_{k,i}(r))}$  is a candidate for parameter  $\nu$ . On the other hand,  $\nu$  needs to satisfy  $\nu \geq \frac{LD}{\sqrt{2}}$  from Theorem 1. Therefore, we set  $\nu$  as follows:

$$\nu = \max \left\{ \sqrt{\sum_{i=1}^N \sum_{r=1}^{n_i} \text{var}(\xi_{k,i}(r))}, \frac{LD}{\sqrt{2}} \right\}.$$

### 5.1.2 Sensitivity analysis

We solve the bandwidth-sharing problem for 12 different settings of parameters shown in Table 1. We consider 4 parameters in our model that scale the problem. Here,  $m_b$  denotes the multiplier of the capacity vector  $b$ ,  $m_c$  denotes the multiplier of the congestion cost function  $c(x)$ , and  $m_\xi$  and  $d_\xi$  are two multipliers that parametrize the random variable  $\xi$ . More precisely, if  $i$ -th user in route  $r$  is uniformly distributed in  $[a - b, a + b]$ , here we assume that it is uniformly distributed in  $[m_\xi a - d_\xi b, m_\xi a + d_\xi b]$ .  $S(i)$  denotes the  $i$ -th setting of parameters. For each of these 4 parameters, we consider 3 settings where one parameter changes and other parameters are fixed. This allows us to observe the sensitivity of the algorithms with respect to each of these parameters. The SA algorithms are terminated after 4000 iterates. To measure the error

-	S(i)	$m_b$	$m_c$	$m_\xi$	$d_\xi$
$m_b$	1	1	1	5	2
	2	0.1	1	5	2
	3	0.01	1	5	2
$m_c$	4	0.1	2	2	1
	5	0.1	1	2	1
	6	0.1	0.5	2	1
$m_\xi$	7	1	1	1	5
	8	1	1	2	5
	9	1	1	5	5
$d_\xi$	10	1	0.01	1	1
	11	1	0.01	1	2
	12	1	0.01	1	5

Table 1: The bandwidth-sharing problem: Parameter settings

of the schemes, we run each scheme 25 times and then compute the mean squared error (MSE) using the metric  $\frac{1}{25} \sum_{i=1}^{25} \|x_k^i - x^*\|^2$  for any  $k = 1, \dots, 4000$ , where  $i$  denotes the  $i$ -th sample. Table 2 shows the 90% confidence intervals (CIs) of the error for the DASA and HSA schemes.

-	S(i)	DASA - 90% CI	HSA with $\theta = 0.1$ - 90% CI	HSA with $\theta = 1$ - 90% CI	HSA with $\theta = 10$ - 90% CI
$m_b$	1	[2.97e-6, 4.66e-6]	[1.52e-6, 2.37e-6]	[1.70e-6, 2.97e-6]	[1.33e-5, 1.81e-5]
	2	[2.97e-6, 4.66e-6]	[1.52e-6, 2.37e-6]	[1.70e-6, 2.97e-6]	[1.33e-5, 1.81e-5]
	3	[1.15e-7, 3.04e-7]	[2.12e-8, 4.92e-8]	[4.66e-8, 1.17e-7]	[8.07e-7, 2.43e-6]
$m_c$	4	[4.39e-7, 6.55e-7]	[1.33e-6, 1.80e-6]	[4.71e-7, 8.75e-7]	[3.84e-6, 5.38e-6]
	5	[1.29e-6, 1.97e-6]	[9.00e-6, 1.20e-5]	[7.88e-7, 1.36e-6]	[5.61e-6, 7.98e-6]
	6	[3.44e-6, 5.36e-6]	[2.26e-4, 2.53e-4]	[1.25e-6, 1.99e-6]	[7.34e-6, 1.12e-5]
$m_\xi$	7	[4.29e-5, 6.40e-5]	[7.92e-5, 1.49e-4]	[2.83e-5, 4.75e-5]	[1.84e-4, 2.75e-4]
	8	[3.18e-5, 4.83e-5]	[3.46e-5, 6.07e-5]	[1.97e-5, 3.39e-5]	[1.40e-4, 1.99e-4]
	9	[1.83e-5, 2.88e-5]	[6.12e-6, 9.99e-6]	[1.06e-5, 1.85e-5]	[8.33e-5, 1.13e-4]
$d_\xi$	10	[3.82e-4, 5.91e-4]	[2.86e+1, 2.86e+1]	[5.50e-1, 5.70e-1]	[7.23e-5, 9.64e-5]
	11	[9.81e-4, 1.44e-3]	[2.86e+1, 2.86e+1]	[5.45e-1, 5.85e-1]	[2.85e-4, 3.80e-4]
	12	[6.26e-3, 8.44e-3]	[2.85e+1, 2.86e+1]	[5.47e-1, 6.44e-1]	[1.77e-3, 2.36e-3]

Table 2: The bandwidth-sharing problem – 90% CIs for DASA and HSA schemes

### 5.1.3 Results and insights

We observe that DASA scheme performs favorably and is far more robust in comparison with the HSA schemes with different choice of  $\theta$ . Importantly, in most of the settings, DASA stands close to the HSA scheme with the minimum MSE. Note that when  $\theta = 1$  or  $\theta = 10$ , the stepsize  $\frac{\theta}{k}$  is not within the interval  $(0, \frac{\eta - \beta L}{(1 + \beta)^2 L^2}]$  for small  $k$  and is not feasible in the sense of Prop. 4. Comparing the performance of each HSA scheme in different settings, we observe that HSA schemes are fairly sensitive to the choice of parameters. For example, HSA with  $\theta = 0.1$  performs very well in settings S(1), S(2), and S(3), while its performance deteriorates in settings S(10), S(11), and S(12). A similar discussion holds for other two HSA schemes. A good instance of this argument is shown in Figure 4. For example, HSA scheme with  $\theta = 10$  performs poorly in settings S(1) and S(4), while it outperforms other schemes in setting S(11). We also observe that changing  $m_b$  from 1 to 0.1 does not affect the error. This is because the optimal solution  $x^*$  remains feasible for a smaller vector  $B$ . On the other hand, the error decreases when we use  $m_b = 0.01$ . Figure 5 presents the flow rates of the users in different routes for the setting S(4). One immediate observation is that the flow rates of HSA scheme with  $\theta = 10$  fluctuates noticeably in the beginning due to a very large stepsize. Figure 6 provides an image of the 90% CIs for the setting S(4). We used two formats to present the intervals. The left-hand side half of each plot shows the intervals with line segments, while the other half shows the lower and upper bound of the intervals continuously. The colorful points represent the 25 sample errors at corresponding iterations. We see that the DASA scheme and HSA scheme with  $\theta = 1$  have CIs with similar size and a smooth mean while the mean in HSA scheme with  $\theta = 10$  is nonsmooth and oscillates more as the algorithm proceeds.



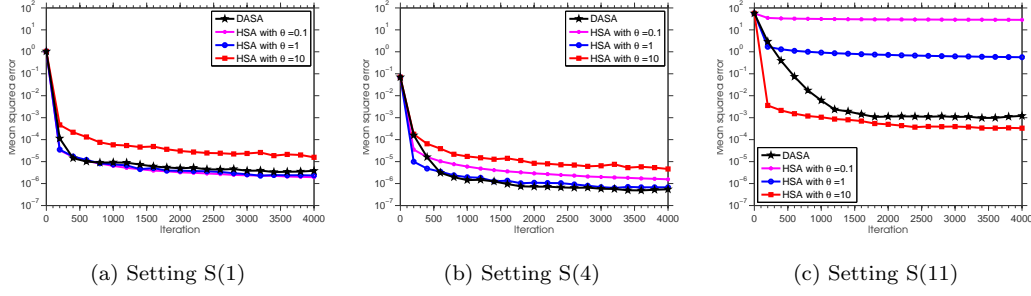


Figure 4: The bandwidth-sharing problem – MSE – DASA vs. HSA schemes

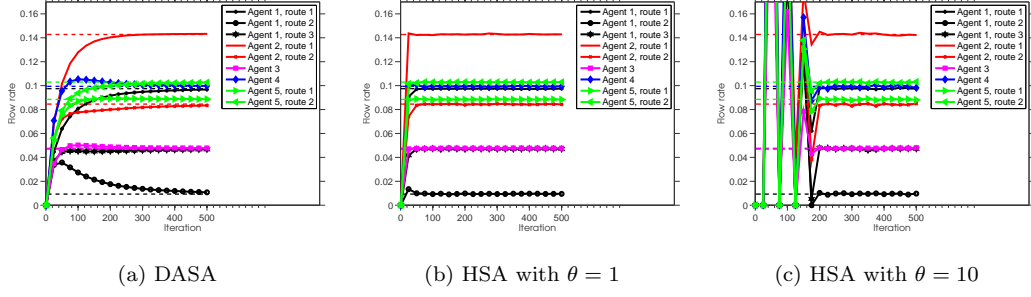


Figure 5: The bandwidth-sharing problem – flow rates for the setting S(4)

## 5.2 A networked stochastic Nash-Cournot game

Consider a networked Nash-Cournot game akin to that described in Example 1. Specifically, let firm  $i$ 's generation and sales decisions at node  $j$  be given by  $g_{ij}$  and  $s_{ij}$ , respectively. Suppose the price function  $p_j$  is given by  $p_j(\bar{s}_j, a_j, b_j) = a_j - b_j \bar{s}_j^\sigma$ , where  $\bar{s}_j = \sum_i s_{ij}$ ,  $\sigma \geq 1$  and  $a_j$  and  $b_j$  are uniformly distributed random variables defined over the intervals  $[lb_j^a, ub_j^a]$  and  $[lb_j^b, ub_j^b]$ , respectively. For purposes of simplicity, we assume that the generation cost is linear and is given by  $c_{ij}g_{ij}$ . We also impose a bound on sales decisions, as specified  $s_{ij} \leq cap_{ij}^l$  for all  $i$  and  $j$ . Note that sales decisions are always bounded by aggregate generation

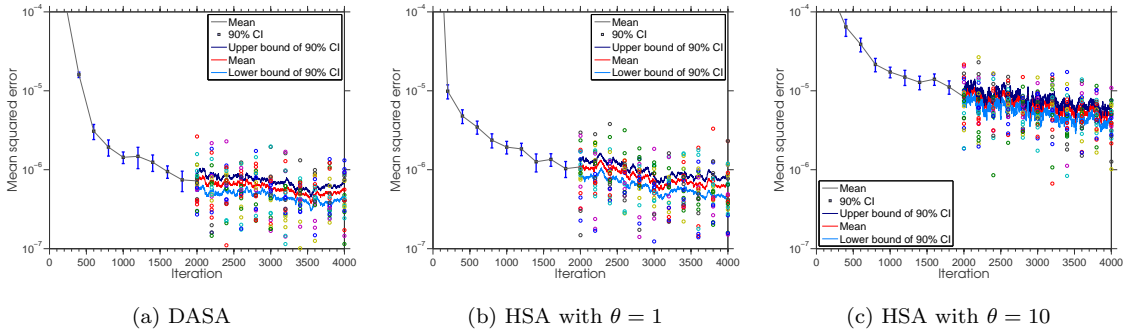


Figure 6: The bandwidth-sharing problem – 90% CIs for the setting S(4)

capacity. The optimization model for the  $i$ -th firm is given by:

$$\begin{aligned} \text{minimize} \quad & \mathbb{E} \left[ \sum_{j=1}^M (c_{ij} g_{ij} - s_{ij} (a_j - b_j \bar{s}_j^\sigma)) \right] \\ \text{subject to} \quad & x_i = (s_i, g_i) \in X_i \triangleq \left\{ \begin{aligned} & \sum_{j=1}^M g_{ij} = \sum_{j=1}^M s_{ij}, \\ & g_{ij} \leq \text{cap}_{ij}, \quad s_{ij} \leq \text{cap}'_{ij}, \quad j = 1, \dots, M, \\ & g_{ij}, s_{ij} \geq 0, \quad j = 1, \dots, M. \end{aligned} \right\}. \end{aligned} \quad (47)$$

As discussed in [15], when  $1 < \sigma \leq 3$  and  $M \leq \frac{3\sigma-1}{\sigma-1}$ , the mapping  $F$  is strictly monotone and strong monotonicity can be induced using a regularized mapping, given that our interest lies in strongly monotone problems. On the other hand, when  $\sigma > 1$ , it is difficult to check that mapping  $F$  has Lipschitzian property. This motivates us to employ the distributed locally randomized SA schemes introduced in Sec. 4.3. Now, using regularization and randomized schemes, we would like to solve the  $\text{VI}(X, F^\epsilon + \eta \mathbf{I})$ , where  $\eta > 0$  is the regularization parameter and  $F^\epsilon$  is defined by (27). As a consequence, this problem admits a unique solution denoted by  $x_{\eta, \epsilon}^*$ .

### 5.2.1 SA algorithms

In this experiment, we use four different SA schemes for solving  $\text{VI}(X, F^\epsilon + \eta \mathbf{I})$  described in Sec. 3.2 and Sec. 4:

**MSR-DASA scheme.** In this scheme, we employ the algorithm (38) and assume that the random vector  $z$  is generated via the MSR scheme, i.e.,  $z_i$  is uniformly distributed on the set  $B_{n_i}(0, \epsilon_i)$  while the mapping  $F_\epsilon$  is defined by (27). One immediate benefit of applying this scheme is that the Lipschitzian parameter can be estimated from Prop. 6b. Moreover, we assume that the stepsizes  $\gamma_{k,i}$  are given by (25)-(26). The multiplier  $r_i$  is randomly chosen for each firm within the prescribed range. The constant  $c$  is maintained at  $\frac{\eta}{4}$ . Parameters  $D$  and  $\nu$  need to be estimated, while the Lipschitzian parameter  $L$  is obtained by Prop. 6b, i.e.,

$$L = \sqrt{N} \|C\| \max_{j=1, \dots, N} \left\{ \kappa_j \frac{n_j!!}{(n_j - 1)!!} \frac{1}{\epsilon_j} \right\}.$$

**MSR-HSA schemes.** Analogous to the MSR-DASA scheme, this scheme uses the distributed locally randomized SA algorithm (38) where for any  $i = 1, \dots, N$ , the random vector  $z_i$  is uniformly drawn from the ball  $B_{n_i}(0, \epsilon_i)$  and mapping  $F_\epsilon$  is defined by (27). The difference is that here we use the harmonic stepsize of the form  $\frac{\theta}{k}$  at  $k$ -th iteration for any firm, where  $\theta > 0$ .

**MCR-DASA scheme.** This scheme is similar to the MSR-DASA scheme with one key difference. We assume that random vector  $z$  is generated by the MCR scheme, i.e., for any  $i = 1, \dots, N$ , random vector  $z_i$  is uniformly drawn from the cube  $C_{n_i}(0, \epsilon_i)$  independent from any  $z_j$  with  $j \neq i$ . The Lipschitz constant  $L$  required for calculating the stepsizes is given by Prop. 7b:

$$L = \frac{\sqrt{n} \|C'\|}{\min_{j=1, \dots, N} \{\epsilon_j\}}.$$

**MCR-HSA schemes.** This scheme uses the algorithm (38) with multi-cubic uniform random variable  $z$ . The stepsizes in this scheme are harmonic of the form  $\frac{\theta}{k}$ .

To obtain the solution  $x_{\eta, \epsilon}^*$ , we use the HSA scheme with the stepsizes  $\frac{1}{k}$  using 20000 iterations. Note that in this experiment, when we use the DASA scheme, we allow that the condition  $\nu \geq \frac{DL}{\sqrt{2}}$  is violated and we replace it with  $\nu \geq D$ . The condition  $\nu \geq D$  keeps the adaptive stepsizes positive for any  $k$ . As a consequence of ignoring  $\nu \geq \frac{DL}{\sqrt{2}}$ , the adaptive stepsizes become larger and in the order of the harmonic stepsizes in our analysis. Note that by this change, the convergence of the DASA algorithm is still guaranteed, while the result of Theorem 1d does not hold necessarily.

### 5.2.2 Sensitivity analysis

We consider a Nash-Cournot game with 5 firms over a network with 3 nodes. We set  $\sigma = 1.1$ ,  $lb_j^b = 0.04$ ,  $ub_j^b = 0.05$ , and  $lb_j^a = 1$  for any  $j$  and  $ub^a = (1.5; 2; 2.5)$ . Having these parameters fixed, our test problems are generated by changing other model's parameters. These parameters are as follows: the parameter of locally randomized schemes  $\epsilon$ , the regularization parameter  $\eta$ , the starting point of the SA algorithm  $x_0$ , and the multiplier  $M_a$  for the random variable  $a_j$  for any  $j$ . We also consider two different settings for  $\text{cap}_{ij}$  and  $\text{cap}'_{ij}$ . Note that when  $\text{cap}_{ij} = 1$ , the constraints  $s_{ij} \leq 3$  are redundant and can be removed. In our analysis we assume that  $\epsilon_i \triangleq \epsilon$  is identical for all firms. Similar to the first experiment in Sec. 5.1.2,

-	S(i)	$\epsilon$	$\eta$	$x_0$	$M_a$	$\text{cap}_{ij}$	$\text{cap}'_{ij}$
$\epsilon$	1	0.1	0.1	$P_1$	1	1	3
	2	0.001	0.1	$P_1$	1	1	3
	3	0.0001	0.1	$P_1$	1	1	3
$\eta$	4	0.1	0.1	$P_2$	1	10	1
	5	0.1	0.05	$P_2$	1	10	1
	6	0.1	0.01	$P_2$	1	10	1
$x_0$	7	0.1	1	$P_1$	6	10	1
	8	0.1	1	$P_2$	6	10	1
	9	0.1	1	$P_3$	6	10	1
$M_a$	10	0.01	0.5	$P_2$	2	1	3
	11	0.01	0.5	$P_2$	4	1	3
	12	0.01	0.5	$P_2$	6	1	3

Table 3: The stochastic Nash-Cournot game – settings of parameters

we consider a set of test problems corresponding to each of these parameters. In each set, one parameter changes and takes 3 different values, while other parameters are fixed. Table 3 represents 12 test problems as described. Note that  $P_1$ ,  $P_2$ , and  $P_3$  are three different feasible starting points. More precisely,  $P_1 = 0$ ,  $P_2 = 0.5(\text{cap}' ; \text{cap})$ , and  $P_3 = (\text{cap}' ; \text{cap})$ . Similar to the first experiment, the termination criteria is running the SA algorithms for 4000 iterates. We run each algorithm 25 times and then we obtain the MSE of the form  $\frac{1}{25} \sum_{i=1}^{25} \|x_k^i - x_{\eta, \epsilon}^*\|^2$  for any  $k = 1, \dots, 4000$ . Table 4 and Table 5 show the 90% CIs of the error for the described schemes.

### 5.2.3 Results and insights

Table 4 presents the simulation results for the test problems using the MSR-DASA and MSR-HSA schemes. One observation is the effect of changing the parameter  $\epsilon$  on the error of the schemes is negligible. We only see a slight change in the error of MSR-HSA scheme with  $\theta = 10$ . Comparing the order of the error, we notice that the MSR-DASA scheme is placed second among all schemes of the first set of the test problems. In the second set, by decreasing  $\eta$  the error of all the schemes, except for the MSR-HSA scheme with  $\theta = 0.1$ , first decreases and then increases. This is not an odd observation since we used  $x_{\eta, \epsilon}^*$  instead of  $x^*$  to measure the errors and  $x_{\eta, \epsilon}^*$  changes itself when  $\eta$  or  $\epsilon$  changes. In this set, the MSR-DASA scheme still has the second best errors among all schemes. The schemes are not much sensitive to the choice of  $x_0$  and we observe that the second place is still reserved by the MSR-DASA scheme. Finally, in the last set, we see that increasing the factor  $M_a$ , as we expect, increases the error in most of the schemes. The reason is that increasing the order of  $M_a$  increases both mean and variance of the random variable  $a$ . Importantly, we observe that our MSR-DASA scheme remains very robust among the MSR-HSA scheme. Table 5 shows the error estimations

-	S(i)	DASA - 90% CI	HSA with $\theta = 0.1$ - 90% CI	HSA with $\theta = 1$ - 90% CI	HSA with $\theta = 10$ - 90% CI
$\epsilon$	1	[1.38e-2, 2.37e-2]	[1.83e+1, 1.87e+1]	[1.60e-1, 2.15e-1]	[3.07e-3, 5.33e-3]
	2	[1.38e-2, 2.37e-2]	[1.83e+1, 1.87e+1]	[1.60e-1, 2.15e-1]	[3.04e-3, 5.30e-3]
	3	[1.38e-2, 2.37e-2]	[1.83e+1, 1.87e+1]	[1.60e-1, 2.15e-1]	[3.04e-3, 5.30e-3]
$\eta$	4	[1.92e-3, 3.98e-3]	[1.63e-0, 1.71e-0]	[8.43e-3, 1.62e-2]	[5.28e-4, 1.08e-3]
	5	[1.42e-3, 3.12e-3]	[1.84e-0, 1.93e-0]	[7.43e-3, 1.44e-2]	[2.59e-4, 5.76e-4]
	6	[5.61e-3, 1.62e-2]	[2.33e-0, 2.44e-0]	[1.61e-2, 2.39e-2]	[5.06e-4, 8.65e-4]
$x_0$	7	[2.68e-6, 3.48e-6]	[4.37e-1, 5.13e-1]	[1.37e-6, 1.92e-6]	[6.71e-6, 9.21e-6]
	8	[2.68e-6, 3.48e-6]	[2.22e-5, 2.91e-5]	[1.37e-6, 1.92e-6]	[6.71e-6, 9.21e-6]
	9	[2.68e-6, 3.48e-6]	[2.22e-5, 2.91e-5]	[1.37e-6, 1.92e-6]	[6.71e-6, 9.21e-6]
$M_a$	10	[4.45e-3, 9.25e-3]	[5.79e-1, 9.25e-1]	[1.67e-3, 5.72e-3]	[2.72e-5, 2.07e-2]
	11	[8.85e-3, 1.73e-2]	[1.25e-0, 2.12e-0]	[9.38e-4, 1.82e-2]	[4.52e-3, 3.22e-2]
	12	[1.92e-2, 3.91e-2]	[8.51e-1, 2.31e-0]	[1.87e-3, 4.15e-2]	[1.04e-2, 7.23e-2]

Table 4: The stochastic Nash-Cournot game – 90% CIs for MSR-DASA and MSR-HSA schemes

using the MCR-DASA and MCR-HSA schemes. Comparing these results with the MSR schemes in Table 4, we see that the sensitivity of the MCR schemes to the parameters is very similar to that of MSR schemes and the MCR-DASA scheme performs as the second best among all MCR schemes. We also see that in most of the settings, the error of the MSR-DASA scheme is slightly smaller than the error of the MCR-DASA

scheme. One reason can be that the MSR scheme has a smaller Lipschitz constant than the MCR scheme for our problem settings.

-	S(i)	DASA - 90% CI	HSA with $\theta = 0.1$ - 90% CI	HSA with $\theta = 1$ - 90% CI	HSA with $\theta = 10$ - 90% CI
$\epsilon$	1	[1.22e-2, 2.55e-2]	[1.84e+1, 1.88e+1]	[1.78e-1, 2.29e-1]	[2.42e-3, 4.21e-3]
	2	[1.21e-2, 2.53e-2]	[1.84e+1, 1.88e+1]	[1.78e-1, 2.28e-1]	[2.37e-3, 4.13e-3]
	3	[1.21e-2, 2.53e-2]	[1.84e+1, 1.88e+1]	[1.78e-1, 2.28e-1]	[2.37e-3, 4.13e-3]
$\eta$	4	[4.17e-3, 9.50e-3]	[1.65e-0, 1.74e-0]	[9.37e-3, 1.84e-2]	[7.38e-4, 1.73e-3]
	5	[1.41e-3, 4.06e-3]	[1.85e-0, 1.93e-0]	[6.88e-3, 1.32e-2]	[2.85e-4, 5.06e-4]
	6	[8.19e-3, 1.88e-2]	[2.37e-0, 2.46e-0]	[1.85e-2, 3.10e-2]	[4.18e-4, 7.05e-4]
$x_0$	7	[2.25e-5, 2.88e-5]	[4.31e-1, 5.12e-1]	[9.41e-6, 1.18e-5]	[3.99e-5, 5.27e-5]
	8	[2.25e-5, 2.88e-5]	[1.13e-4, 1.58e-4]	[9.40e-6, 1.18e-6]	[3.99e-5, 5.27e-5]
	9	[2.25e-5, 2.88e-5]	[1.13e-4, 1.58e-4]	[9.40e-6, 1.18e-5]	[3.99e-5, 5.27e-5]
$M_a$	10	[1.66e-3, 4.29e-3]	[6.17e-1, 8.88e-1]	[4.21e-4, 1.79e-3]	[3.82e-4, 8.30e-3]
	11	[3.03e-3, 1.22e-2]	[1.29e-0, 2.23e-0]	[9.63e-4, 5.77e-3]	[2.48e-3, 2.52e-2]
	12	[6.05e-3, 3.26e-2]	[8.50e-1, 2.49e-0]	[2.27e-3, 1.29e-2]	[5.54e-3, 5.67e-2]

Table 5: The stochastic Nash-Cournot game – 90% CIs for MCR-DASA and MCR-HSA schemes

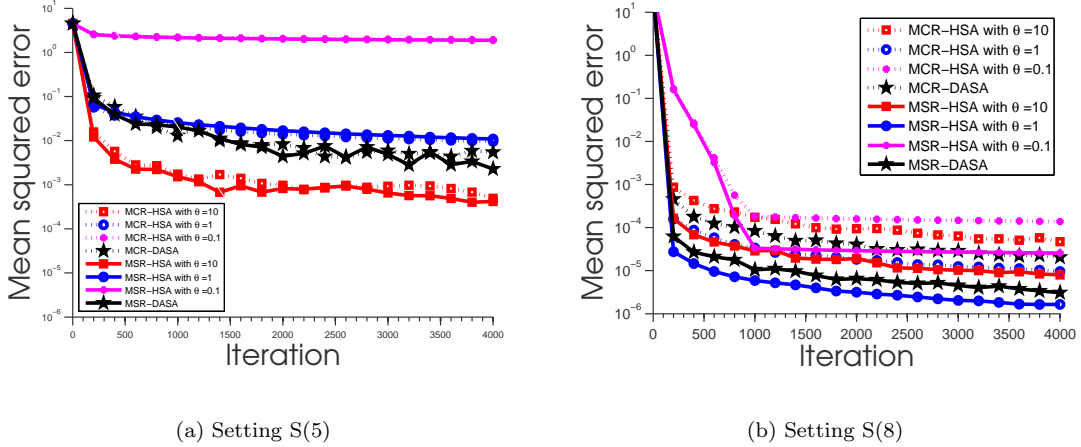


Figure 7: The stochastic Nash-Cournot game – comparison among all the schemes

Figure 7 illustrates a comparison among the different schemes described in Sec. 5.2.1 for the case of setting S(5) and S(8). All the MSR schemes are shown with solid lines, while the MCR schemes are presented with dashed lines. There are some immediate observations here. Regarding the order of the error, in both if the settings S(5) and S(8), the schemes with the distributed adaptive stepsizes given by (25)-(26) are the second best scheme among each of MSR and MCR schemes. This indicates the robustness of the DASA scheme compared with the HSA schemes. We also observe that in the setting S(5), the HSA schemes with  $\theta = 10$  (both MSR and MCR) have the minimum error, while in setting S(8), the HSA schemes with  $\theta = 1$  has the minimum error. This is an illustration of sensitivity of HSA schemes to the setting of problem parameters. Let us now compare the MSR schemes with the MCR schemes. In the setting S(5), the MSR and MCR schemes perform very closely and in fact, it is hard to distinguish the difference between their errors. On the other hand, in the setting S(8), we see that the MSR schemes have a better performance than their MCR counterparts. Figure 8 illustrates the 90% confidence intervals for the MSR schemes with the setting S(5). Teese intervals are shown with line segments in the left-hand side half of each plot and shown with continious bouns in the right-hand side half. The colourful points present the samples at each level of iterates. Impostantly, we observe that the CIs of MSR-DASA scheme are as tight as the MSR-HSA scheme with  $\theta = 1$  and they are tighter than the ones in the MSR-HSA scheme with  $\theta = 10$ . Figure 9 shows the similar comparison for the MCR schemes.

## 6 Concluding remarks

We consider the solution of strongly monotone Cartesian stochastic variational inequality problems through stochastic approximation (SA) schemes. Motivated by the naive stepsize rules employed in most SA implementations, we develop a recursive rule that adapts to problem parameters such as the Lipschitz and

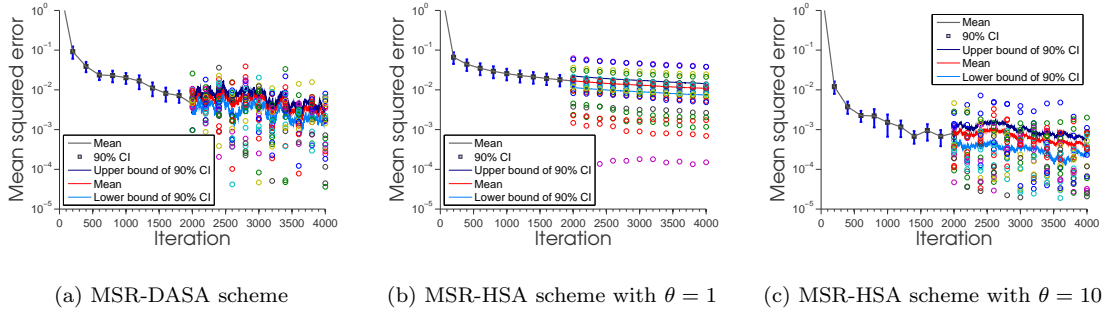


Figure 8: The stochastic Nash-Cournot game – setting S(5) – MSR-DASA vs. MSR-HSA schemes

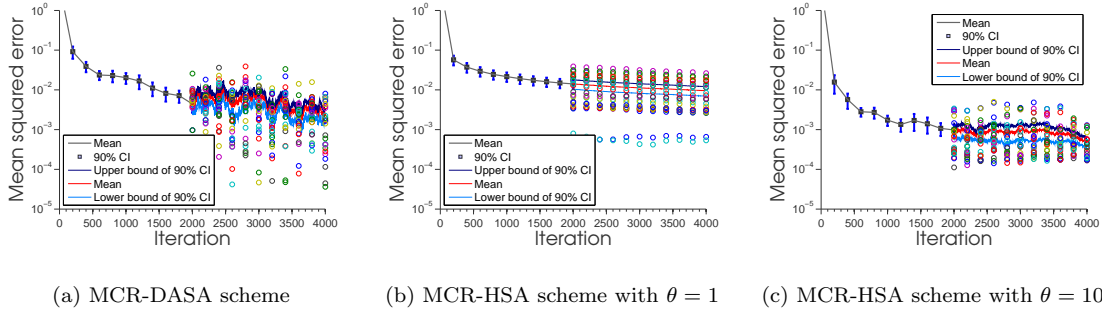


Figure 9: The stochastic Nash-Cournot game – setting S(5) – MCR-DASA vs. MCR-HSA schemes

monotonicity constants of the map and ensures almost-sure convergence of the iterates to the unique solution. An extension to the distributed multi-agent regime is provided. A shortcoming of this approach is the reliance on the availability of a Lipschitz constant. This motivates the construction of two locally randomized techniques to cope with instances where the mapping is either not Lipschitz or estimating the parameter is challenging. In each of these techniques, we show that an approximation of the original mapping is Lipschitz continuous with a prescribed constant. We utilize these techniques in developing a distributed locally randomized adaptive steplength SA scheme where we perturb the mapping at each iteration by a uniform random variable over a prescribed distribution. It is shown that this scheme produces iterates that converge to a solution of an approximate problem, and the sequence of approximate solutions converge to the unique solution of the original stochastic variational problem. In Sec. 5, we apply our schemes on two sets of problems, a bandwidth-sharing problem in communication networks and a networked stochastic Nash-Cournot game. Through these examples, we observed that the adaptive distributed stepsize scheme displays far more robustness than the standard implementations that leverage harmonic stepsizes of the form  $\frac{\theta}{k}$  in both problems. Furthermore, the randomized smoothing techniques assume utility in the Cournot regime where Lipschitz constants cannot be easily derived.

## References

- [1] T. ALPCAN AND T. BAŞAR, *A game-theoretic framework for congestion control in general topology networks*, in Proceedings of the 41st IEEE Conference on Decision and Control, December 2002, pp. 1218–1224.
- [2] —, *Distributed algorithms for Nash equilibria of flow control games*, in Advances in Dynamic Games, vol. 7 of Annals of the International Society of Dynamic Games, Birkhäuser Boston, 2003, pp. 473–498.

- [3] D. P. BERTSEKAS, *Stochastic optimization problems with nondifferentiable functionals with an application in stochastic programming*, in Proceedings of 1972 IEEE Conference on Decision and Control, 1972, pp. 555–559.
- [4] V. S. BORKAR, *Stochastic Approximation: A Dynamical Systems Viewpoint*, Cambridge University Press, 2008.
- [5] R. H. BYRD, M. E. HRIBAR, AND J. NOCEDAL, *An interior point algorithm for large-scale nonlinear programming*, SIAM Journal on Optimization, 9 (1999), pp. 877–900.
- [6] S.-W. CHO AND A. GOEL, *Bandwidth allocation in networks: a single dual update subroutine for multiple objectives*, Combinatorial and algorithmic aspects of networking, 3405 (2005), pp. 28–41.
- [7] D. CICEK, M. BROADIE, AND A. ZEEVI, *General bounds and finite-time performance improvement for the kiefer-wolfowitz stochastic approximation algorithm*, To appear in Operations Research, (2011).
- [8] Y. M. ERMOLIEV, *Stochastic Programming Methods*, Nauka, Moscow, 1976.
- [9] ———, *Stochastic quasigradient methods*, in Numerical Techniques for Stochastic Optimization, Springer-Verlag, 1983, pp. 141–185.
- [10] ———, *Stochastic quasigradient methods and their application to system optimization*, Stochastics, 9 (1983), pp. 1–36.
- [11] F. FACCHINEI AND J.-S. PANG, *Finite-dimensional variational inequalities and complementarity problems. Vols. I,II*, Springer Series in Operations Research, Springer-Verlag, New York, 2003.
- [12] A. P. GEORGE AND W. B. POWELL, *Adaptive stepsizes for recursive estimation with applications in approximate dynamic programming*, Machine Learning, 65 (2006), pp. 167–198.
- [13] A. M. GUPAL, *Stochastic methods for solving nonsmooth extremal problems (Russian)*, Naukova Dumka, 1979.
- [14] H. JIANG AND H. XU, *Stochastic approximation approaches to the stochastic variational inequality problem*, IEEE Transactions on Automatic Control, 53 (2008), pp. 1462–1475.
- [15] A. KANNAN AND U. V. SHANBHAG, *Distributed computation of equilibria in monotone Nash games via iterative regularization techniques*, SIAM Journal of Optimization, 22 (2012), pp. 1177–1205.
- [16] A. KANNAN, U. V. SHANBHAG, AND H. . M. KIM, *Addressing supply-side risk in uncertain power markets: stochastic Nash models, scalable algorithms and error analysis*, Optimization Methods and Software (online first), 0 (2012), pp. 1–44.
- [17] A. KANNAN, U. V. SHANBHAG, AND H. M. KIM, *Strategic behavior in power markets under uncertainty*, Energy Systems, 2 (2011), pp. 115–141.
- [18] F. KELLY, A. MAULLOO, AND D. TAN, *Rate control for communication networks: shadow prices, proportional fairness, and stability*, Journal of the Operational Research Society, 49 (1998), pp. 237–252.
- [19] J. KOSHAL, A. NEDIĆ, AND U. V. SHANBHAG, *Single timescale regularized stochastic approximation schemes for monotone nash games under uncertainty*, Proceedings of the IEEE Conference on Decision and Control (CDC), (2010), pp. 231–236.
- [20] J. KOSHAL, A. NEDIĆ, AND U. V. SHANBHAG, *Single timescale stochastic approximation for stochastic nash games in cognitive radio systems*, 2011, pp. 1–8.
- [21] H. J. KUSHNER AND G. G. YIN, *Stochastic Approximation and Recursive Algorithms and Applications*, Springer New York, 2003.
- [22] H. LAKSHMANAN AND D. FARIAS, *Decentralized resource allocation in dynamic networks of agents*, SIAM Journal on Optimization, 19 (2008), pp. 911–940.

- [23] J. LINDEROTH, A. SHAPIRO, AND S. WRIGHT, *The empirical behavior of sampling methods for stochastic programming*, Ann. Oper. Res., 142 (2006), pp. 215–241.
- [24] C. METZLER, B. HOBBS, AND J.-S. PANG, *Nash-cournot equilibria in power markets on a linearized dc network with arbitrage: Formulations and properties*, Networks and Spatial Theory, 3 (2003), pp. 123–150.
- [25] V. I. NORKIN, *The analysis and optimization of probability functions*, tech. rep., International Institute for Applied Systems Analysis technical report, 1993. WP-93-6.
- [26] B. POLYAK, *Introduction to optimization*, Optimization Software, Inc., New York, 1987.
- [27] W. B. POWELL, *Approximate Dynamic Programming: Solving the Curses of Dimensionality*, John Wiley & Sons, Inc., Hoboken, New Jersey, 2010.
- [28] H. ROBBINS AND S. MONRO, *A stochastic approximation method*, Ann. Math. Statistics, 22 (1951), pp. 400–407.
- [29] G. SCUTARI, F. FACCHINEI, J.-S. PANG, AND D. PALOMAR, *Monotone games for cognitive radio systems*, (2012), pp. 83–112.
- [30] G. SCUTARI AND D. P. PALOMAR, *Mimo cognitive radio: a game theoretical approach*, IEEE Transactions on Signal Processing, 58 (2010), pp. 761–780.
- [31] S. SHAKKOTTAI AND R. SRIKANT, *Network optimization and control*, Foundations and Trends in Networking, 2 (2007), pp. 271–379.
- [32] U. V. SHANBHAG, G. INFANGER, AND P. GLYNN, *A complementarity framework for forward contracting under uncertainty*, Operations Research, 59 (2011), pp. 810–834.
- [33] A. SHAPIRO, *Monte Carlo sampling methods*, in Handbook in Operations Research and Management Science, vol. 10, Elsevier Science, Amsterdam, 2003, pp. 353–426.
- [34] J. C. SPALL, *Introduction to Stochastic Search and Optimization: Estimation, Simulation, and Control*, Wiley, Hoboken, NJ, 2003.
- [35] R. SRIKANT, *Mathematics of Internet Congestion Control*, Birkhauser, 2004.
- [36] V. A. STEKLOV, *Sur les expressions asymptotiques decertaines fonctions dfinies par les quations differentielles du second ordre et leers applications au problme du dveloppement d’une fonction arbitraire en sries procant suivant les diverses fonctions*, Comm. Charkov Math. Soc., 2 (1907), pp. 97–199.
- [37] ———, *Main Problems of Mathematical Physics*, Nauka, Moscow, 1983.
- [38] J. WANG, G. SCUTARI, AND D. P. PALOMAR, *Robust mimo cognitive radio via game theory*, IEEE Transactions on Signal Processing, 59 (2011), pp. 1183–1201.
- [39] H. YIN, U. V. SHANBHAG, AND P. G. MEHTA, *Nash equilibrium problems with scaled congestion costs and shared constraints*, IEEE Transactions of Automatic Control, 56 (2009), pp. 1702–1708.
- [40] F. YOUSEFIAN, A. NEDIĆ, AND U. SHANBHAG, *A regularized adaptive steplength stochastic approximation scheme for monotone stochastic variational inequalities*, Proceedings of the 2011 Winter Simulation Conference, (2011), pp. 4110–4121.
- [41] F. YOUSEFIAN, A. NEDIĆ, AND U. V. SHANBHAG, *On stochastic gradient and subgradient methods with adaptive steplength sequences*, Automatica, 48 (2012), pp. 56–67. An extended version of the paper available at: <http://arxiv.org/abs/1105.4549>.